

AD-A065 285

STANFORD UNIV CALIF DEPT OF COMPUTER SCIENCE
SPARSE AND PARALLEL MATRIX COMPUTATIONS.(U)
DEC 78 F T LUK
STAN-CS-78-685

F/G 12/1

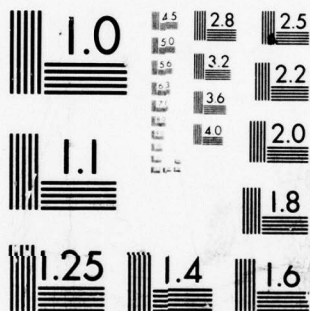
DAHC04-75-6-0185
NL

UNCLASSIFIED

1 OF 2

AD
A065 285





MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

LEVEL

13
NW

SPARSE AND PARALLEL MATRIX COMPUTATIONS

by

Franklin Tai-cheung Luk

AD A0 65285

DDC FILE COPY

STAN-CS-78-685
DECEMBER 1978



COMPUTER SCIENCE DEPARTMENT
School of Humanities and Sciences
STANFORD UNIVERSITY

This document has been approved
for public release and sale; its
distribution is unlimited.



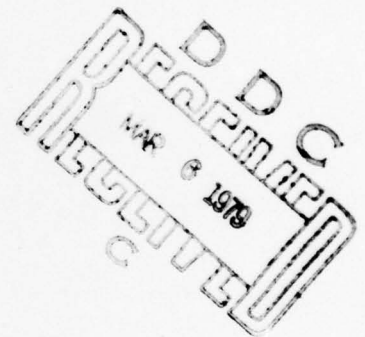
79 3 01 124

SPARSE AND PARALLEL MATRIX COMPUTATIONS

by

Franklin Tai-cheung Luk

ACCESSION NO.	
DTIS	White Section <input checked="" type="checkbox"/>
DOC	Buff Section <input type="checkbox"/>
UNANNOUNCED	<input type="checkbox"/>
JUSTIFICATION	
BY	
DISTRIBUTION/AVAILABILITY STATE	
50	
A	



This document has been approved
for public release and sale; its
distribution is unlimited.

Research supported in part under U. S. Army Research Office grant
DAHCO4-75-G-0185.

79 03 01 12 4

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER STAN-CS-78-685	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) SPARSE AND PARALLEL MATRIX COMPUTATIONS		5. TYPE OF REPORT & PERIOD COVERED Technical, December 1978
7. AUTHOR(s) Franklin Tai-cheung Luk		6. PERFORMING ORG. REPORT NUMBER STAN-CS-78-685
9. PERFORMING ORGANIZATION NAME AND ADDRESS Computer Science Department Stanford University Stanford, California 94305		8. CONTRACT OR GRANT NUMBER(s) DAHCO4-75-G-0185
11. CONTROLLING OFFICE NAME AND ADDRESS U.S. Army Research Office Box 12211 Triangle Park, N.C. 27709		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		12. REPORT DATE December 1978
		13. NUMBER OF PAGES 168 p.
		15. SECURITY CLASS. (of this report) Unclassified
16. DISTRIBUTION STATEMENT (of this Report) Releasable without limitations on dissemination.		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) This thesis deals with four important matrix problems: (1) the application of many variants of the conjugate gradient method for solving matrix equations, (2) the solution of lower and upper bounds quadratic programs associated with M-matrices, (3) the construction of a Block Lanczos method for computing the greatest singular values of a matrix, and (4) the computation of the singular value decomposition of a matrix on the ILLIAC-IV computer.		

DD FORM 1 JAN 73 1473

EDITION OF 1 NOV 65 IS OBSOLETE

UNCLASSIFIED 094 120 RLH

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

ABSTRACT.

This thesis deals with four important matrix problems: (1) the application of many variants of the conjugate gradient method for solving matrix equations, (2) the solution of lower and upper bounds quadratic programs associated with M-matrices, (3) the construction of a Block Lanczos method for computing the greatest singular values of a matrix, and (4) the computation of the singular value decomposition of a matrix on the ILLIAC-IV computer.

ACKNOWLEDGMENTS

I am deeply indebted to Professor Gene Golub for teaching me about numerical linear algebra and for providing guidance and support throughout my years of graduate study.

I appreciate the suggestions of Professors Joseph Oliger and James Wilkinson which have greatly improved my presentation.

Three researchers have aided me in essential ways. The many discussions with Professor Marcello Pagano led to the work on quadratic programming. Dr. Louis Hageman and Professor David Young shared with me their knowledge of iterative methods.

I would like to thank all the Serra House residents, particularly Mr. William Coughran and Mrs. Janet Wright, for their help and friendship.

The expert typing and editing skills of Miss Rosemarie Stampfel have made my task of writing much easier.

With sincere gratitude I wish to thank my parents, for starting me on my studies in this country; Gene and Rosemarie, for making life so pleasant; and my wife Vivian, for her love.

TABLE OF CONTENTS

	PAGE
I. INTRODUCTION	1
II. CONJUGATE GRADIENT ACCELERATION OF ITERATIVE METHODS	3
2.1. Introduction	3
2.2. Conjugate Gradient Method	8
2.3. Preconditioned Conjugate Gradient Method...	16
2.4. Equivalence Results	22
2.5. Block Jacobi Methods	28
2.6. A Compressed Method and Cyclic Reduction	32
2.7. Computer Implementation	41
2.8. Applications of the RS-CG Method	50
References	59
III. QUADRATIC PROGRAMMING WITH M-MATRICES	60
3.1. Introduction	60
3.2. Linear Complementarity Problem	64
3.3. Lower and Upper Bounds Problem	76
3.4. Problem with Non-finite Bounds	82
3.5. Numerical Examples	85
3.6. Conclusion	94
References	95
IV. A BLOCK LANCZOS METHOD FOR COMPUTING SINGULAR VALUES	96
4.1. Introduction	96
4.2. Block Lanczos Method for Symmetric Matrices	97
4.3. Block Bidiagonalization Algorithm	102
4.4. Error Bounds for the Singular Value Approximations	108
4.5. Iterating to Improve Accuracy	114
4.6. Block Bidiagonalization Method with Reorthogonalization	118
4.7. Iterative Block Lanczos Method	122
4.8. Test Examples	125
References	132

TABLE OF CONTENTS Continued

	PAGE
V. COMPUTING THE SINGULAR VALUE DECOMPOSITION ON THE ILLIAC IV	134
5.1. Introduction	134
5.2. The ILLIAC IV Computer	135
5.3. Programming Languages for the ILLIAC.....	136
5.4. A Row Orthogonalization Method	137
5.5. Least Squares Solutions	146
5.6. Data Structures	151
5.7. Numerical Properties	153
5.8. Test Results	154
References	158

TABLE OF ALGORITHMS

	PAGE
II. CONJUGATE GRADIENT ACCELERATION OF ITERATIVE METHODS	
Algorithm 2.1. CG method	10
Algorithm 3.1.	17
Algorithm 3.2. Preconditioned CG method or CG-accelerated single method.....	18
Algorithm 4.1. CG-accelerated double method.....	23
Algorithm 4.2. Derived CG-accelerated double method.	25
Algorithm 6.1. CJ-CG method	33
Algorithm 6.2. RS-CG method	37
Algorithm 7.1. CG method	41
Algorithm 7.2. Preconditioned CG method	43
Algorithm 7.3. Generalized CG method	45
Algorithm 7.4. RS-CG method	47
III. QUADRATIC PROGRAMMING WITH M-MATRICES	
Algorithm 2.1. Chandrasekaran's method	68
Algorithm 2.2. Modified Chandrasekaran's method....	74
Algorithm 2.3.	75
Algorithm 3.1. Pang's method	78
Algorithm 3.2. Modified Pang's method	81
IV. A BLOCK LANCZOS METHOD FOR COMPUTING SINGULAR VALUES	
Algorithm 2.1. Block Lanczos method	98
Algorithm 3.1. Block Bidiagonalization method.....	106
Algorithm 5.1.	117
Algorithm 6.1. Block Bidiagonalization method with Reorthogonalization	121
Algorithm 7.1. Iterative Block Lanczos method.....	124
V. COMPUTING THE SINGULAR VALUE DECOMPOSITION ON THE ILLIAC-IV	
Algorithm 1. SVD	142
Algorithm 2. MINFIT	148

TABLE OF ABBREVIATIONS

<u>Abbreviation</u>	<u>Name</u>	<u>First appears on page</u>
CG method	Conjugate Gradient method	1
J-CG method	CG-accelerated single block Jacobi method	29
DJ-CG method	CG-accelerated double block Jacobi method	29
CJ-CG method	Compressed J-CG method	33
RS-CG method	CG method applied to reduced system of linear equations	36
SVD	Singular Value Decomposition	134
MINFIT	Algorithm for Least Squares problems	148

I. INTRODUCTION

This thesis consists of five chapters. In the next four chapters, we discuss the solution of four important matrix problems under the assumption that they are all of large order. For the first three problems we construct algorithms which exploit the sparsity of the associated matrices. We use a computer with parallel processing abilities to solve the fourth problem.

Chapter 2 of this thesis deals with the many variants of the conjugate gradient (CG) method for solving matrix equations. First, we describe a new preconditioned CG method and derive its many attractive properties. Second, we introduce the CG-accelerated single and double methods and give conditions for their mathematical equivalence. Third, we apply the equivalence result to systems of linear equations possessing 'Property A'. Fourth, this original result is used to show that a new method based on conjugate gradients and cyclic reduction is equivalent to a popular method of Reid. Fifth, we demonstrate that our new method is more efficient than Reid's method in both work and storage.

In Chapter 3 we study the problem of quadratic programming with M-matrices. We describe (1) an effective algorithm for the case where the variables are subject to a lower bound constraint, and (2) an analogous algorithm for the case where the variables are subject to lower and upper bounds constraints. We demonstrate the special monotone behavior of the iterate and gradient vectors.

The result on the gradient vector is new. It leads us to consider a simple updating procedure which preserves the monotonicity of both vectors. The procedure uses the fact that an M-matrix has a non-negative inverse. Two new algorithms are then constructed by incorporating this updating procedure into the two given algorithms. We give numerical examples which show that the new methods can be twice as fast as the original ones.

A Block Lanczos method is introduced in Chapter 4 for computing a few greatest singular values and associated vectors of a matrix. We present a theoretical development of the method and give a theorem on its rate of convergence. The practical implementation aspects are then discussed and particular attention is paid to the choice of block size. We believe that all our results are original. We are in fact unaware of any other algorithms which solve this problem.

In the fifth chapter, we study the computation of the singular value decomposition of a matrix on the ILLIAC IV computer. We describe the architecture of the machine and explain why the standard Golub-Reinsch algorithm is not applicable to this problem. We then present a one-sided orthogonalization method which makes very efficient use of the parallel computing abilities of the ILLIAC machine. Our method is shown to be Jacobi-like and numerically stable. Finally, a comparison of our method on the ILLIAC IV computer with the Golub-Reinsch algorithm on a conventional machine demonstrates the great potential of parallel computers in the important area of matrix computations.

II. CONJUGATE GRADIENT ACCELERATION OF ITERATIVE METHODS

2.1. Introduction

In this chapter we are concerned with the application of the conjugate gradient method (CG method) as an iterative solution procedure (cf. Reid [8]) for large and sparse sets of linear equations.

Let us consider the system

$$A\tilde{x} = \tilde{b} , \quad (1.1)$$

where A is a given $n \times n$ nonsingular matrix and \tilde{b} a given column vector. We assume that the matrix A is large and sparse. Instead of (1.1) we may solve the modified system

$$C^{-1}A\tilde{x} = C^{-1}\tilde{b} , \quad (1.2)$$

where C is some $n \times n$ nonsingular matrix such that it is a simple computing task to solve the system

$$C\tilde{z} = \tilde{d} .$$

We usually choose C so that the new system (1.2) is "better-conditioned" with respect to inversion than the original system (1.1) (cf. [1]), i.e.

$$\kappa(C^{-1}A) < \kappa(A) ,$$

where $\kappa(M) = \|M\| \cdot \|M^{-1}\|$ for any nonsingular matrix M . We refer to $\kappa(M)$ as the condition number of M with respect to inversion (and with respect to the matrix norm $\|\cdot\|$). The matrix C is frequently

an approximation of A so that the condition number of $C^{-1}A$ is small, e.g., C can be a product of sparse triangular matrices (see [7]).

We wish to apply the conjugate gradient method to the solution of the preconditioned system (1.2). We refer to the new method as the preconditioned conjugate gradient method (preconditioned CG method). This problem has been studied by many researchers (see [1] for a bibliography), who considered the case where A and C are both symmetric and positive definite matrices. We make a slightly weaker assumption here.

ASSUMPTION 1. There exists some nonsingular matrix S of order n such that the matrix $SC^{-1}AS^{-1}$ is symmetric and positive definite.

We call S a "symmetrization" matrix. The next lemma shows that Assumption 1 is meaningful.

LEMMA 1.1. Suppose that both matrices A and C are symmetric and positive definite. If we choose the matrix S such that $S^t S = C$, then Assumption 1 is valid.

Proof. Let $B \equiv SC^{-1}AS^{-1}$. The matrix B is symmetric because

$$B = S(S^t S)^{-1}AS^{-1} = (S^{-1})^t AS^{-1},$$

and is positive definite because

$$y^t B y = (S^{-1}y)^t A(S^{-1}y) > 0 \quad \text{for } y \neq 0. \quad \square$$

This apparently new concept of a "symmetrization" matrix makes it a straightforward exercise (in Section 2.3) to derive the properties of the preconditioned CG method from those of the CG method. We see in the same section that the matrix S need not be formed explicitly and that only the product $S^t S$ is required for the preconditioned CG method. For many problems, several choices of the "symmetrization" matrix are possible (see, e.g., Section 2.4). An interesting problem is therefore to determine the effect of the "symmetrization" matrix on the convergence rate of the preconditioned CG method.

Letting

$$G \equiv I - C^{-1}A \quad (1.3)$$

and

$$\underline{k} \equiv C^{-1}\underline{b}, \quad (1.4)$$

we can rewrite the system (1.2) as

$$(I - G)\underline{x} = \underline{k}, \quad (1.5)$$

from which we derive the iterative scheme

$$\underline{x}^{(i+1)} = G\underline{x}^{(i)} + \underline{k} \quad \text{for } i = 0, 1, \dots, \quad (1.6)$$

where $\underline{x}^{(0)}$ is some initial approximation to \underline{x} . We refer to (1.6) as the single method for the system (1.5). Let us write the matrix A as

$$A = D + L + U,$$

where D is a (block) diagonal matrix, L is a strictly lower triangular matrix and U is a strictly upper triangular matrix. We then call (1.6) the (block) Jacobi method if $C = D$, the (block) Gauss-Seidel method if $C = D + L$, or the symmetric (block) successive over-relaxation method (SSOR method) if

$$C = (D + \omega L) D^{-1} (D + \omega U) ,$$

where ω is a scalar parameter. Other choices of C lead to other well-known single methods.

We can combine the iterates from the single method to define a more general iterative procedure

$$\tilde{x}_w^{(i)} = \sum_{j=0}^i v_j^{(i)} \tilde{x}^{(j)} \quad \text{for } i = 0, 1, \dots, \quad (1.7)$$

where the $v_j^{(i)}$'s are constants such that

$$\sum_{j=0}^i v_j^{(i)} = 1 \quad \text{for } i = 0, 1, \dots . \quad (1.8)$$

We call (1.7) a semi-iterative method with respect to the method (1.6) (cf. [11]). A particular instance of (1.7) is the second-order procedure

$$\tilde{x}^{(i+1)} = \alpha_{i+1} [\beta_i (\tilde{G} \tilde{x}^{(i)} + \tilde{k}) + (1 - \beta_i) \tilde{x}^{(i)}] + (1 - \alpha_{i+1}) \tilde{x}^{(i-1)} , \quad (1.9)$$

for $i = 0, 1, \dots ,$

where $\alpha_1 = 1$. Many methods can be expressed in the form of (1.9), e.g. the Chebyshev semi-iterative method and the Richardson second-order method (see [4]). In the next section, we show that we can also express the conjugate gradient method in this three-term recurrence form. We refer to (1.9) as the CG-accelerated single method if the iteration parameters $\{\alpha_i\}$ and $\{\beta_i\}$ are chosen based on the CG method.

Combining two iterations of the single method (1.6), we obtain the double method

$$x^{(i+1)} = G^2 x^{(i)} + G\tilde{k} + \tilde{k} \quad \text{for } i = 0, 1, \dots, \quad (1.10)$$

where $x^{(0)} = \tilde{x}^{(0)}$. The iterative procedure (1.10) can be accelerated by the second-order procedure

$$x^{(i+1)} = r_{i+1} [\rho_i (G^2 x^{(i)} + G\tilde{k} + \tilde{k}) + (1-\rho_i) x^{(i)}] + (1-r_{i+1}) x^{(i-1)} \quad (1.11)$$

for $i = 0, 1, \dots$,

where $r_1 = 1$. We want to use the CG method to compute the iteration parameters $\{r_i\}$ and $\{\rho_i\}$, in which case the procedure (1.11) is referred to as the CG-accelerated double method. We need the following assumption.

ASSUMPTION 2. There exists some nonsingular matrix T of order n such that the matrix $T(I - G^2)T^{-1}$ is symmetric and positive definite.

It can be easily checked that $T = S$ is one possible choice.

We have five goals to accomplish in this chapter: (1) to describe the preconditioned CG method and to derive its many nice properties, (2) to determine conditions under which the CG-accelerated single method and the CG-accelerated double method are "virtually-equivalent," i.e.

$$\tilde{x}^{(2i)} = x^{(i)} \quad \text{for } i = 0, 1, \dots, \quad (1.12)$$

(3) to apply the equivalence result of (2) to the case of the block Jacobi method when the coefficient matrix A is symmetric and positive definite, and can be written in partitioned form as

$$A = \left(\begin{array}{c|c} M_1 & F \\ \hline F^t & M_2 \end{array} \right) ,$$

where M_1 and M_2 are square matrices such that the systems

$$M_1 \tilde{z}_1 = \tilde{d}_1 \quad \text{and} \quad M_2 \tilde{z}_2 = \tilde{d}_2$$

are easy to solve, (4) to introduce a new method based on conjugate gradients and cyclic reduction for the class of matrix equations considered in (3), and to use the result of (3) to show that our new method is mathematically equivalent to a generalization of a method due to Reid [9], and (5) to demonstrate that our new method is more efficient than the generalized Reid's method in both work and storage, and thus is an effective solution procedure for an important class of matrix equations. We believe that all five results are original.

2.2. Conjugate Gradient Method

Let us consider the system

$$A\tilde{x} = \tilde{b} . \quad (1.1)$$

We assume that the matrix A is symmetric and positive definite, and that the matrix C has been chosen as the identity matrix. Our goal is to use the conjugate gradient method to compute the iteration

parameters $\{\alpha_i\}$ and $\{\beta_i\}$ in the second-order procedure

$$\tilde{x}^{(i+1)} = \alpha_{i+1} [\beta_i (G\tilde{x}^{(i)} + \tilde{k}) + (1-\beta_i)\tilde{x}^{(i)}] + (1-\alpha_{i+1})\tilde{x}^{(i-1)}, \quad (1.9)$$

for $i = 0, 1, \dots$,

where $\alpha_1 = 1$. We define the residual vector $\tilde{r}^{(i)}$ by

$$\tilde{r}^{(i)} \equiv \tilde{b} - A\tilde{x}^{(i)} \quad \text{for } i = 0, 1, \dots. \quad (2.1)$$

As we have

$$G = I - A \quad \text{and} \quad \tilde{k} = \tilde{b},$$

it follows that

$$G\tilde{x}^{(i)} + \tilde{k} = \tilde{x}^{(i)} + \tilde{r}^{(i)},$$

and the procedure (1.9) becomes

$$\tilde{x}^{(i+1)} = \alpha_{i+1} (\beta_i \tilde{r}^{(i)} + \tilde{x}^{(i)}) + (1 - \alpha_{i+1})\tilde{x}^{(i-1)} \quad (2.2)$$

for $i = 0, 1, \dots$,

where $\alpha_1 = 1$. Using the notation

$$(p, q) \equiv p^t q,$$

we can express the conjugate gradient method for solving the system (1.1) as follows (cf. Rutishauser [10]).

ALGORITHM 2.1 (CG method)

(1) Let $\tilde{x}^{(0)}$ be an initial approximation to x . Compute

$$\begin{aligned}\tilde{r}^{(0)} &:= b - A\tilde{x}^{(0)}, \\ \beta_0 &:= \frac{(\tilde{r}^{(0)}, \tilde{r}^{(0)})}{(\tilde{r}^{(0)}, A\tilde{r}^{(0)})}\end{aligned}$$

and

$$\tilde{x}^{(1)} := \beta_0 \tilde{r}^{(0)} + \tilde{x}^{(0)}.$$

Let $\alpha_1 := 1$ and $i := 0$.

(2) Repeat until $\tilde{r}^{(i)} = 0$:

(a) Let

$$i := i + 1.$$

(b) Compute

$$\begin{aligned}\tilde{r}^{(i)} &:= b - A\tilde{x}^{(i)}, \\ \beta_i &:= \frac{(\tilde{r}^{(i)}, \tilde{r}^{(i)})}{(\tilde{r}^{(i)}, A\tilde{r}^{(i)})}\end{aligned}$$

and

$$\alpha_{i+1} := \left[1 - \frac{\beta_i}{\beta_{i-1}} \frac{(\tilde{r}^{(i)}, \tilde{r}^{(i)})}{(\tilde{r}^{(i-1)}, \tilde{r}^{(i-1)})} \frac{1}{\alpha_i} \right]^{-1}.$$

(c) Compute

$$\tilde{x}^{(i+1)} := \alpha_{i+1} [\beta_i \tilde{r}^{(i)} + \tilde{x}^{(i)}] + (1 - \alpha_{i+1}) \tilde{x}^{(i-1)}. \quad \square$$

The residual vectors satisfy the well-known orthogonality property (see [6])

$$(\tilde{r}^{(i)}, \tilde{r}^{(j)}) = 0 \quad \text{for } i \neq j. \quad (2.3)$$

Thus, in the absence of round-off errors the solution vector is obtained in at most n iterations of the CG-method. But the computed residuals are not orthogonal in practice. Our approach is to permit the gradual loss of orthogonality and with it the finite termination property of the method. We are concerned primarily with the iterative aspects of the CG algorithm. In fact, when used for solving large sparse sets of linear equations arising from the discretization of elliptic partial differential equations, the CG method often converges within a number of iterations small compared with n (see, e.g., [2]). Nonetheless, this orthogonality property (2.3) is of theoretical importance for it characterizes the CG method among all second-order procedures of the form (2.2).

THEOREM 2.1. The second-order iterative procedure

$$\tilde{x}^{(i+1)} = \alpha_{i+1}(\beta_i \tilde{r}^{(i)} + \tilde{x}^{(i)}) + (1 - \alpha_{i+1})\tilde{x}^{(i-1)} \quad \text{for } i = 0, 1, \dots,$$

where $\alpha_1 = 1$ and

$$\tilde{r}^{(i)} = \tilde{b} - A\tilde{x}^{(i)} \quad \text{for } i = 0, 1, \dots,$$

reduces to the conjugate gradient method (Algorithm 2.1) if and only if

$$(\tilde{r}^{(i)}, \tilde{r}^{(j)}) = 0 \quad \text{for } i \neq j.$$

Proof. We only have to prove the sufficiency part. Let $i = 0$.

From (2.2), we obtain

$$\tilde{x}^{(1)} = \beta_0 \tilde{r}^{(0)} + \tilde{x}^{(0)}.$$

Thus,

$$\underline{b} - A\underline{x}^{(1)} = -\beta_0 A\underline{r}^{(0)} + \underline{b} - A\underline{x}^{(0)},$$

which becomes

$$\underline{x}^{(1)} = \underline{x}^{(0)} - \beta_0 A\underline{r}^{(0)}. \quad (2.4)$$

As $(\underline{x}^{(0)}, \underline{x}^{(1)}) = 0$, we find

$$0 = (\underline{x}^{(0)}, \underline{x}^{(0)}) - \beta_0 (\underline{x}^{(0)}, A\underline{r}^{(0)}).$$

It follows that

$$\beta_0 = \frac{(\underline{x}^{(0)}, \underline{x}^{(0)})}{(\underline{x}^{(0)}, A\underline{r}^{(0)})}.$$

Now let $i \geq 1$. From (2.2), we get

$$\underline{x}^{(i+1)} = \alpha_{i+1}(-\beta_i A\underline{r}^{(i)} + \underline{x}^{(i)}) + (1 - \alpha_{i+1})\underline{x}^{(i-1)}. \quad (2.5)$$

Therefore,

$$0 = \alpha_{i+1}[-\beta_i (\underline{x}^{(i)}, A\underline{r}^{(i)}) + (\underline{x}^{(i)}, \underline{x}^{(i)})]$$

because

$$(\underline{x}^{(i)}, \underline{x}^{(i+1)}) = (\underline{x}^{(i)}, \underline{x}^{(i-1)}) = 0.$$

Consequently,

$$\beta_i = \frac{(\underline{x}^{(i)}, \underline{x}^{(i)})}{(\underline{x}^{(i)}, A\underline{r}^{(i)})}.$$

From (2.5), as $(\underline{x}^{(i-1)}, \underline{x}^{(i+1)}) = (\underline{x}^{(i-1)}, \underline{x}^{(i)}) = 0$, we get

$$0 = -\alpha_{i+1}\beta_i (\underline{x}^{(i-1)}, A\underline{r}^{(i)}) + (1 - \alpha_{i+1})(\underline{x}^{(i-1)}, \underline{x}^{(i-1)}). \quad (2.6)$$

Since $\alpha_1 = 1$, we can replace i by $i-1$ in (2.5) to obtain

$$\tilde{x}^{(i)} = \alpha_i (-\beta_{i-1} A \tilde{x}^{(i-1)} + \tilde{x}^{(i-1)}) + (1 - \alpha_i) \tilde{x}^{(i-2)},$$

or, equivalently,

$$A \tilde{x}^{(i-1)} = -\frac{1}{\beta_{i-1}} \left\{ \frac{1}{\alpha_i} [\tilde{x}^{(i)} - (1 - \alpha_i) \tilde{x}^{(i-2)}] - \tilde{x}^{(i-1)} \right\}.$$

Therefore,

$$\begin{aligned} (\tilde{x}^{(i-1)}, A \tilde{x}^{(i)}) &= (A \tilde{x}^{(i-1)}, \tilde{x}^{(i)}) \\ &= -\frac{1}{\alpha_i \beta_{i-1}} (\tilde{x}^{(i)}, \tilde{x}^{(i)}) \end{aligned}$$

because $(\tilde{x}^{(i-2)}, \tilde{x}^{(i)}) = (\tilde{x}^{(i-1)}, \tilde{x}^{(i)}) = 0$. We can now simplify (2.6) to

$$0 = \frac{\alpha_{i+1} \beta_i}{\alpha_i \beta_{i-1}} (\tilde{x}^{(i)}, \tilde{x}^{(i)}) + (1 - \alpha_{i+1}) (\tilde{x}^{(i-1)}, \tilde{x}^{(i-1)}),$$

or equivalently,

$$0 = \left[\frac{\beta_i}{\beta_{i-1}} \frac{(\tilde{x}^{(i)}, \tilde{x}^{(i)})}{(\tilde{x}^{(i-1)}, \tilde{x}^{(i-1)})} \frac{1}{\alpha_i} - 1 \right] \alpha_{i+1} + 1.$$

We conclude that

$$\alpha_{i+1} = \left[1 - \frac{\beta_i}{\beta_{i-1}} \frac{(\tilde{x}^{(i)}, \tilde{x}^{(i)})}{(\tilde{x}^{(i-1)}, \tilde{x}^{(i-1)})} \frac{1}{\alpha_i} \right]^{-1}. \quad \square$$

From the recurrence relations (2.4) and (2.5), we can verify by induction that

$$\tilde{x}^{(i+1)} = \tilde{x}^{(0)} - A \left[\sum_{j=0}^i \sigma_j^{(i+1)} A^j \right] \tilde{x}^{(0)} \quad \text{for } i = 0, 1, \dots, \quad (2.7)$$

where the $\sigma_j^{(i+1)}$'s are constants. As

$$\tilde{x}^{(i+1)} - \tilde{x}^{(0)} = -A(\tilde{x}^{(i+1)} - \tilde{x}^{(0)}),$$

it follows that

$$\tilde{x}^{(i+1)} = \tilde{x}^{(0)} + \left[\sum_{j=0}^i \sigma_j^{(i+1)} A^j \right] \tilde{x}^{(0)} \quad \text{for } i = 0, 1, \dots \quad (2.8)$$

The conjugate gradient method is therefore a special instance of the polynomial acceleration procedure

$$\tilde{x}^{(i+1)} = \tilde{x}^{(0)} + P_i(A) \tilde{x}^{(0)} \quad \text{for } i = 0, 1, \dots, \quad (2.9)$$

where $P_i(A)$ is a matrix polynomial in A of degree i . Let us define an error function

$$f_M(\chi) \equiv (\chi - \tilde{x}, M(\chi - \tilde{x})), \quad (2.10)$$

where M is some $n \times n$ positive definite matrix. The conjugate gradient method possesses the following optimality property (see [3, p. 397]).

THEOREM 2.2. Among all polynomial acceleration procedures of the form (2.9), the conjugate gradient method generates an iterate $\tilde{x}^{(i)}$ that minimizes the error function $f_A(\tilde{x}^{(i)})$ of (2.10), for $i = 1, 2, \dots$.

A few interesting results follow from Theorem 2.2. Let us assume that the matrix A has only $p < n$ distinct eigenvalues. Then there exists a matrix polynomial $Q_p(A)$ of degree p so that

$$Q_p(A) = 0.$$

Therefore,

$$f_A(\tilde{x}^{(p)}) = 0$$

and the CG method converges in at most p steps. The same result also holds if the matrix A has a large number of distinct eigenvalues but the initial error vector

$$\tilde{e}^{(0)} = \tilde{x}^{(0)} - \tilde{x}$$

lies in a subspace generated by the eigenvectors associated with only p of these eigenvalues. The CG method is thus optimal for the particular eigenvector mix of the initial error $\tilde{e}^{(0)}$. As the iteration proceeds, the extreme eigenvalues are approximated especially well and the CG procedure would then behave as if the corresponding eigenvectors were not present. Hence we often observe a superlinear convergence rate for the CG method (see, e.g., [2]). It follows that for rapid convergence the eigenvalues should be sparse at the extremes and dense in the center. Convergence would be slow if the eigenvalues were packed in the extremes.

Another consequence of Theorem 2.2 is an error estimate of the CG method. Assume that the eigenvalues of matrix A are included in the interval $[\alpha, \beta]$, where $\alpha > 0$. Then by choosing an appropriate Chebyshev polynomial acceleration, we get (see [3, p. 428])

$$f_A(\tilde{x}^{(i)}) \leq \frac{1}{T_i^2\left(\frac{\beta + \alpha}{\beta - \alpha}\right)} f_A(\tilde{x}^{(0)}), \quad (2.11)$$

where T_i is the i -th degree Chebyshev polynomial of the first kind.

Since

$$\begin{aligned} 2 \cdot T_i \left(\frac{\beta + \alpha}{\beta - \alpha} \right) &= \left(\frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^i + \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^i \\ &\geq \left(\frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^i, \end{aligned}$$

where $\kappa = \beta/\alpha$, the inequality (2.11) becomes

$$f_A(\tilde{x}^{(i)}) \leq 4 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2i} f_A(\tilde{x}^{(0)}) . \quad (2.12)$$

For fixed i , the error bound of (2.12) decreases as κ decreases, and is thus tightest when α and β are the extreme eigenvalues of A , in which case the number κ is the condition number of A with respect to the spectral norm. The bound is, however, pessimistic asymptotically.

The above convergence properties of the CG method are important when we consider the choice of a preconditioning matrix in the next section.

2.3. Preconditioned Conjugate Gradient Method

We wish to apply the conjugate gradient method to solve the preconditioned system

$$C^{-1}Ax = C^{-1}b . \quad (1.2)$$

We construct the matrix equation

$$SC^{-1}AS^{-1}y = SC^{-1}b , \quad (3.1)$$

where $y = Sx$ and S is the "symmetrization" matrix of Assumption 1.

Now we apply the conjugate gradient method to the system (3.1).

ALGORITHM 3.1.

(1) Let $\underline{u}^{(0)}$ be an initial approximation to \underline{u} . Compute

$$\begin{aligned}\tilde{\underline{r}}^{(0)} &:= SC^{-1}\underline{b} - SC^{-1}AS^{-1}\underline{u}^{(0)} , \\ \beta_0 &:= \frac{(\tilde{\underline{r}}^{(0)}, \tilde{\underline{r}}^{(0)})}{(\tilde{\underline{r}}^{(0)}, A\tilde{\underline{r}}^{(0)})}\end{aligned}$$

and

$$\underline{u}^{(1)} := \beta_0 \tilde{\underline{r}}^{(0)} + \underline{u}^{(0)} .$$

Let $\alpha_1 := 1$ and $i := 0$.

(2) Repeat until $\tilde{\underline{r}}^{(i)} = \underline{0}$:

(a) Let

$$i := i + 1 .$$

(b) Compute

$$\begin{aligned}\tilde{\underline{r}}^{(i)} &:= SC^{-1}\underline{b} - SC^{-1}AS^{-1}\underline{u}^{(i)} , \\ \beta_i &:= \frac{(\tilde{\underline{r}}^{(i)}, \tilde{\underline{r}}^{(i)})}{(\tilde{\underline{r}}^{(i)}, SC^{-1}AS^{-1}\tilde{\underline{r}}^{(i)})}\end{aligned}$$

and

$$\alpha_{i+1} := \left[1 - \frac{\beta_i}{\beta_{i-1}} \frac{(\tilde{\underline{r}}^{(i)}, \tilde{\underline{r}}^{(i)})}{(\tilde{\underline{r}}^{(i-1)}, \tilde{\underline{r}}^{(i-1)})} \frac{1}{\alpha_i} \right]^{-1} .$$

(c) Compute

$$\underline{u}^{(i+1)} := \alpha_{i+1} (\beta_i \tilde{\underline{r}}^{(i)} + \underline{u}^{(i)}) + (1 - \alpha_{i+1}) \underline{u}^{(i-1)} . \quad \square$$

Let us define the iterate vector $\underline{x}^{(i)}$ by

$$\underline{x}^{(i)} \equiv S^{-1}\underline{u}^{(i)} \quad \text{for } i = 0, 1, \dots , \quad (3.2)$$

and the pseudo-residual vector $\underline{z}^{(i)}$ by

$$\tilde{z}^{(i)} \equiv C^{-1}\tilde{b} - C^{-1}A\tilde{x}^{(i)} \quad \text{for } i = 0, 1, \dots \quad (3.3)$$

Substituting $S\tilde{x}^{(i)}$ for $\tilde{u}^{(i)}$ and $S\tilde{z}^{(i)}$ for $\tilde{r}^{(i)}$ in Algorithm 3.1, we get an iterative procedure for

$$C^{-1}A\tilde{x} = C^{-1}\tilde{b} \quad (1.2)$$

ALGORITHM 3.2. (Preconditioned CG method or CG-accelerated single method)

(1) Let $\tilde{x}^{(0)}$ be an initial approximation to \tilde{x} . Solve the system

$$C\tilde{z}^{(0)} = \tilde{b} - A\tilde{x}^{(0)}.$$

Compute

$$\beta_0 := \frac{(S\tilde{z}^{(0)}, S\tilde{z}^{(0)})}{(S\tilde{z}^{(0)}, SC^{-1}A\tilde{z}^{(0)})}$$

and

$$\tilde{x}^{(1)} := \beta_0\tilde{z}^{(0)} + \tilde{x}^{(0)}.$$

Let $\alpha_1 := 1$ and $i := 0$.

(2) Repeat until $\tilde{z}^{(i)} = 0$:

(a) Let $i := i + 1$.

(b) Solve the system

$$C\tilde{z}^{(i)} = \tilde{b} - A\tilde{x}^{(i)}.$$

(c) Compute

$$\beta_i := \frac{(S\tilde{z}^{(i)}, S\tilde{z}^{(i)})}{(S\tilde{z}^{(i)}, SC^{-1}A\tilde{z}^{(i)})},$$

$$\alpha_{i+1} := \left[1 - \frac{\beta_i}{\beta_{i-1}} \frac{(S\tilde{z}^{(i)}, S\tilde{z}^{(i)})}{(S\tilde{z}^{(i-1)}, S\tilde{z}^{(i-1)})} \frac{1}{\alpha_i} \right]^{-1}$$

and

$$\tilde{x}^{(i+1)} := \alpha_{i+1}(\beta_i \tilde{z}^{(i)} + \tilde{x}^{(i)}) + (1 - \alpha_{i+1})\tilde{x}^{(i-1)}. \quad \square$$

Observe that we need not form the matrix S explicitly. It suffices to have the product $S^t S$. Also, note that the preconditioned CG method reduces to the generalized CG method of Concus et al. [2] if A and C are symmetric and positive definite matrices, and if $S^t S = C$.

From (2.3) and (3.3), we get that the pseudo-residual vectors satisfy the conjugacy property

$$(S\tilde{z}^{(i)}, S\tilde{z}^{(j)}) = 0 \quad \text{for } i \neq j. \quad (3.4)$$

Since the matrix $S^t S$ is symmetric and positive definite, the preconditioned CG method therefore terminates in at most n steps in exact arithmetic. In practice, the calculated pseudo-residuals do not satisfy (3.4) due to round-off errors. Our approach is to permit the gradual loss of conjugacy and to consider primarily the iterative aspect of the preconditioned CG method.

The next theorem, a direct consequence of Theorem 2.1, states how the conjugacy property of (3.4) characterizes the preconditioned CG method.

THEOREM 3.1. The second-order iterative procedure

$$\tilde{x}^{(i+1)} = \alpha_{i+1}(\beta_1 \tilde{z}^{(i)} + \tilde{x}^{(i)}) + (1 - \alpha_{i+1})\tilde{x}^{(i-1)} \quad \text{for } i = 0, 1, \dots,$$

where $\alpha_1 = 1$ and

$$\tilde{z}^{(i)} = C^{-1}(b - A\tilde{x}^{(i)}) \quad \text{for } i = 0, 1, \dots,$$

reduces to the preconditioned conjugate gradient method (Algorithm 3.2)

if and only if

$$(S\tilde{z}^{(i)}, S\tilde{z}^{(j)}) = 0 \quad \text{for } i \neq j. \quad \square$$

From the relation (2.8) we get

$$\tilde{u}^{(i+1)} = \tilde{u}^{(0)} + \left[\sum_{j=0}^i \tau_j^{(i+1)} H^j \right] \tilde{r}^{(0)} \quad \text{for } i = 0, 1, \dots,$$

where $H \equiv SC^{-1}AS^{-1}$ and the $\tau_j^{(i+1)}$'s are constants. Thus,

$$S\tilde{x}^{(i+1)} = S\tilde{x}^{(0)} + S \left[\sum_{j=0}^i \tau_j^{(i+1)} K^j \right] S^{-1} S\tilde{z}^{(0)} \quad \text{for } i = 0, 1, \dots,$$

where $K \equiv C^{-1}A$, so that

$$\tilde{x}^{(i+1)} = \tilde{x}^{(0)} + \left[\sum_{j=0}^i \tau_j^{(i+1)} K^j \right] \tilde{z}^{(0)} \quad \text{for } i = 0, 1, \dots \quad (3.5)$$

The preconditioned CG method is therefore a special instance of the

polynomial acceleration procedure

$$\tilde{x}^{(i+1)} = \tilde{x}^{(0)} + Q_i(K) \tilde{z}^{(0)}, \quad (3.6)$$

where $Q_i(K)$ is a matrix polynomial in $K (\equiv C^{-1}A)$ of degree i . Recall the error function

$$f_M(\underline{y}) \equiv (\underline{y} - \underline{x}, M(\underline{y} - \underline{x})) , \quad (2.10)$$

where M is a positive definite matrix. As

$$(\underline{u}^{(i)} - \underline{u}, SC^{-1}AS^{-1}(\underline{u}^{(i)} - \underline{u})) = (\underline{x}^{(i)} - \underline{x}, S^tSC^{-1}A(\underline{x}^{(i)} - \underline{x})),$$

the following optimality property of the precondition CG method can be derived from Theorem 2.2.

THEOREM 3.2. Among all polynomial acceleration procedures of the form (3.6), the preconditioned CG method generates an iterate $\underline{x}^{(i)}$ that minimizes the error function $f_B(\underline{x}^{(i)})$ of (2.10), where $B \equiv S^tSC^{-1}A$, for $i = 1, 2, \dots$.

We can draw from Theorem 3.2 conclusions similar to those from Theorem 2.2. The assumption that the matrix $SC^{-1}AS^{-1}$ is symmetric and positive definite implies that the eigenvalues of $C^{-1}A$ are all real and positive. Had the matrix $C^{-1}A$ only $p < n$ distinct eigenvalues, the preconditioned CG method would converge in at most p steps. The eigenvalue distribution of $C^{-1}A$ determines the rate at which the method converges. Convergence would be fast if the eigenvalues were sparse in the extremes and dense in the interior. We often observe that the method converges superlinearly.

Using arguments similar to those in [3, p. 428] and the last section, we obtain an error estimate for the preconditioned CG method:

$$f_B(\tilde{x}^{(i)}) \leq 4 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2i} f_B(\tilde{x}^{(0)}) , \quad (3.7)$$

where $B \equiv S^t S C^{-1} A$ and κ is the spectral condition number of $C^{-1}A$.

A good preconditioning matrix C should therefore have the following features:

- (a) it is a simple computing task to solve matrix equations with coefficient matrix C ,
- (b) the matrix $C^{-1}A$ has a favorable eigenvalue distribution, i.e. either that the eigenvalues are sparse in the extremes and dense in the interior or that only a few of them are distinct, and
- (c) the spectral condition number of $C^{-1}A$ is much smaller than that of A .

2.4. Equivalence Results

Consider the double method

$$\tilde{y}^{(i+1)} = G^2 \tilde{y}^{(i)} + G\tilde{k} + \tilde{k} \quad \text{for } i = 0, 1, \dots , \quad (1.10)$$

where $\tilde{y}^{(0)} = \tilde{x}^{(0)}$. Recall that

$$G \equiv I - C^{-1}A \quad (1.3)$$

and

$$\tilde{k} \equiv C^{-1}b . \quad (1.4)$$

The procedure (1.10) can be regarded as the single method for the system

$$(I - G^2)\underline{x} = G\underline{k} + \underline{k} . \quad (4.1)$$

Suppose that we have chosen a nonsingular matrix T such that the matrix $T(I - G^2)T^{-1}$ is symmetric and positive definite. Let us apply the conjugate gradient method to the system

$$[T(I - G^2)T^{-1}]T\underline{x} = T(G\underline{k} + \underline{k}) .$$

ALGORITHM 4.1. (CG-accelerated double method)

(1) Let

$$\underline{y}^{(0)} := \underline{x}^{(0)} .$$

Compute

$$\underline{w}^{(0)} := G\underline{k} + \underline{k} - (I - G^2)\underline{y}^{(0)} ,$$

$$\rho_0 := \frac{(T\underline{w}^{(0)}, T\underline{w}^{(0)})}{(T\underline{w}^{(0)}, T(I - G^2)\underline{w}^{(0)})}$$

and

$$\underline{y}^{(1)} := \rho_0 \underline{w}^{(0)} + \underline{y}^{(0)} .$$

Let $\gamma_1 := 1$ and $i := 0$.

(2) Repeat until $\underline{w}^{(i)} = 0$:

(a) Let

$$i := i + 1 .$$

(b) Compute

$$\underline{w}^{(i)} := G\underline{k} + \underline{k} - (I - G^2)\underline{y}^{(i)} ,$$

(c) Compute

$$\rho_i := \frac{(T_{\tilde{W}}^{(i)}, T_{\tilde{W}}^{(i)})}{(T_{\tilde{W}}^{(i)}, T(I - G^2) \tilde{W}^{(i)})} ,$$

$$r_{i+1} := \left[1 - \frac{\rho_i}{\rho_{i-1}} \frac{(T_{\tilde{W}}^{(i)}, T_{\tilde{W}}^{(i)})}{(T_{\tilde{W}}^{(i-1)}, T_{\tilde{W}}^{(i-1)})} \frac{1}{r_i} \right]^{-1}$$

and

$$x^{(i+1)} := r_{i+1} [\rho_i \tilde{w}^{(i)} + x^{(i)}] + (1 - r_{i+1}) x^{(i-1)} . \quad \square$$

We now assume that the condition

$$\beta_i = 1 \quad \text{for } i = 0, 1, \dots , \quad (4.2)$$

holds for the CG-accelerated single method (Algorithm 3.2). The three-term recurrence defining the method thereby simplifies to

$$\tilde{x}^{(i+1)} = \alpha_{i+1} (G\tilde{x}^{(i)} + \tilde{k}) + (1 - \alpha_{i+1}) \tilde{x}^{(i-1)} \quad \text{for } i = 0, 1, \dots , \quad (4.3)$$

where $\alpha_1 = 1$. Replacing i by $2j + 1$ in (4.3), we get

$$\tilde{x}^{(2j+2)} = \alpha_{2j+2} (G\tilde{x}^{(2j+1)} + \tilde{k}) + (1 - \alpha_{2j+2}) \tilde{x}^{(2j)} .$$

Also from (4.3), we obtain

$$G\tilde{x}^{(2j+1)} = \alpha_{2j+1} (G^2 \tilde{x}^{(2j)} + G\tilde{k}) + (1 - \alpha_{2j+1}) G\tilde{x}^{(2j-1)}$$

and

$$G\tilde{x}^{(2j-1)} = \frac{1}{\alpha_{2j}} \tilde{x}^{(2j)} - \left(\frac{1 - \alpha_{2j}}{\alpha_{2j}} \right) \tilde{x}^{(2j-2)} - \tilde{k} .$$

It follows that

$$\begin{aligned} \tilde{x}^{(2j+2)} &= \alpha_{2j+2} \alpha_{2j+1} (G^2 \tilde{x}^{(2j)} + G\tilde{k} + \tilde{k}) \\ &\quad + \left[1 + \frac{\alpha_{2j+2}}{\alpha_{2j}} (1 - \alpha_{2j+1}) - \alpha_{2j+2} \right] \tilde{x}^{(2j)} \\ &\quad - \frac{\alpha_{2j+2}}{\alpha_{2j}} (1 - \alpha_{2j+1}) (1 - \alpha_{2j}) \tilde{x}^{(2j-2)} . \end{aligned} \quad (4.4)$$

We use (4.4) to define one more iterative scheme.

ALGORITHM 4.2. (Derived CG-accelerated double method)

(1) Let

$$\tilde{y}^{(0)} := \tilde{x}^{(0)} .$$

Compute

$$\begin{aligned} \hat{\tilde{w}}^{(0)} &:= G\tilde{k} + \tilde{k} - (I - G^2) \tilde{y}^{(0)} , \\ \hat{\rho}_0 &:= \alpha_2 \alpha_1 \end{aligned}$$

and

$$\tilde{y}^{(1)} := \hat{\rho}_0 \hat{\tilde{w}}^{(0)} + \tilde{y}^{(0)} .$$

Let

$$i := 0 .$$

(2) Repeat until $\hat{\tilde{w}}^{(i)} = \underline{0}$:

(a) Let

$$i := i + 1 .$$

(b) Compute

$$\hat{\tilde{w}}^{(i)} := G\tilde{k} + \tilde{k} - (I - G^2) \tilde{y}^{(i)} .$$

(c) Compute

$$\hat{\gamma}_{i+1} := 1 + \frac{\alpha_{2i+2}}{\alpha_{2i}} (1 - \alpha_{2i+1})(1 - \alpha_{2i}) ,$$

$$\hat{\rho}_i := \frac{\alpha_{2i+2} \alpha_{2i+1}}{\hat{\gamma}_{i+1}}$$

and

$$\underline{y}^{(i+1)} := \hat{\gamma}_{i+1} [\hat{\rho}_i \hat{\underline{w}}^{(i)} + \underline{y}^{(i)}] + (1 - \hat{\gamma}_{i+1}) \underline{y}^{(i-1)} . \quad \square$$

We observe that Algorithms 4.1 and 4.2 are identical in form. Our goal is to prove that they are mathematically equivalent, i.e.

$$\underline{x}^{(i)} = \underline{y}^{(i)} \quad \text{for } i = 0, 1, \dots .$$

Because of Theorem 3.1, we only have to show that

$$(T_{\hat{\underline{w}}}^{(i)}, T_{\hat{\underline{w}}}^{(j)}) = 0 \quad \text{for } i \neq j ,$$

where

$$\hat{\underline{w}}^{(i)} = G\underline{k} + \underline{k} - (I - G^2) \underline{y}^{(i)} \quad \text{for } i = 0, 1, \dots .$$

Let us define the spectral radius of an $n \times n$ matrix M by the number

$$\rho(M) \equiv \max_{1 \leq i \leq n} |\lambda_i| ,$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of M . We make the assumption:

ASSUMPTION 3.

$$\rho(G) < 1 .$$

The matrix $(I + G)$ is therefore nonsingular because all its eigenvalues are bounded away from zero. We now choose the "symmetrization" matrix T of Assumption 2 as

$$T = S(I + G)^{-1}, \quad (4.5)$$

where S is the "symmetrization" matrix of Assumption 1. It remains to verify that the matrix $B \equiv T(I - G^2)T^{-1}$ is symmetric and positive definite. As

$$T(I - G^2)T^{-1} = I - (SGS^{-1})^2,$$

the matrix B is symmetric. Hence the eigenvalues of B are real. By Assumption 3, the eigenvalues of (SGS^{-1}) are all less than unity in modulus. Therefore, all the eigenvalues of B are positive.

Since the derived CG-accelerated double method has been constructed directly from the CG-accelerated single method, we have that

$$\tilde{y}^{(i)} = \tilde{x}^{(2i)} \quad \text{for } i = 0, 1, \dots$$

Now,

$$\begin{aligned} \hat{\tilde{w}}^{(i)} &= G\tilde{k} + \tilde{k} - (I - G^2) \tilde{x}^{(2i)} \\ &= (I + G) \tilde{z}^{(2i)} \quad \text{for } i = 0, 1, \dots, \end{aligned}$$

where

$$\tilde{z}^{(2i)} \equiv \tilde{k} - C^{-1}A\tilde{x}^{(2i)} \quad \text{for } i = 0, 1, \dots \quad (3.3)$$

Thus,

$$T\hat{\tilde{w}}^{(i)} = S\tilde{z}^{(2i)}.$$

But from Theorem 3.1, we see that

$$(S\tilde{z}^{(2i)}, S\tilde{z}^{(2j)}) = 0 \quad \text{for } i \neq j.$$

We have thereby proved the desired equivalence result.

THEOREM 4.1. The CG-accelerated single method and the CG-accelerated double method are "virtually-equivalent" if the iteration parameters $\{\beta_i\}$ of the single method satisfy

$$\beta_i = 1 \quad \text{for } i = 0, 1, \dots,$$

and if the "symmetrization" matrices S and T of Assumptions 1 and 2 are chosen so that

$$T = S(I + G)^{-1}.$$

Note that another possibility for the "symmetrization" matrix T is

$$T = S,$$

but then the above "virtual-equivalence" result no longer holds.

2.5. Block Jacobi Methods

In this and the ensuing section, we consider symmetric and positive definite matrices of the form

$$A = \left(\begin{array}{c|c} M_1 & F \\ \hline F^t & M_2 \end{array} \right), \quad (5.1)$$

where M_1 and M_2 are square matrices such that the systems

$$M_1 z_1 = d_1 \quad \text{and} \quad M_2 z_2 = d_2$$

are easy to solve. Many matrices can be described by (5.1), e.g. those matrices that possess "Property A" [11].

The block Jacobi method corresponds to the preconditioning matrix

$$C = \left(\begin{array}{c|c} M_1 & 0 \\ \hline 0 & M_2 \end{array} \right) . \quad (5.2)$$

As both matrices M_1 and M_2 are symmetric and positive definite, it follows that C is, too. We have

$$G = \left(\begin{array}{c|c} 0 & -M_1^{-1}F \\ \hline -M_2^{-1}F^t & 0 \end{array} \right) \quad (5.3)$$

and

$$\tilde{k} = \left(\begin{array}{c} M_1^{-1}b_1 \\ M_2^{-1}b_2 \end{array} \right) . \quad (5.4)$$

Our goal is to show that the CG-accelerated single block Jacobi method (J-CG method) and the CG-accelerated double block Jacobi method (DJ-CG method) are "virtually-equivalent."

There exist matrices S_1 and S_2 such that

$$S_1^t S_1 = M_1 \quad \text{and} \quad S_2^t S_2 = M_2 . \quad (5.5)$$

We choose the "symmetrization" matrix S as

$$S = \left(\begin{array}{c|c} S_1 & 0 \\ \hline 0 & S_2 \end{array} \right) \quad (5.6)$$

All is well because $S^t S = C$ (see Lemma 1.1). Also, let

$$B \equiv SC^{-1}AS^{-1}.$$

Since

$$B = \left(\begin{array}{c|c} I & (S_1^{-1})^t F S_2^{-1} \\ \hline (S_2^{-1})^t F^t S_1^{-1} & I \end{array} \right),$$

its eigenvalues are either unity or of the form $1 \pm \sigma^2$, where σ is some nonzero real number. As the matrix B is positive definite, it follows that $\sigma^2 < 1$. Now,

$$G = I - SBS^{-1},$$

and so the eigenvalues of G are either zero or of the form $\pm \sigma^2$.

Hence

$$\rho(G) < 1,$$

i.e. Assumption 3 is valid. We can therefore choose the matrix $S(I + G)^{-1}$ as the "symmetrization" matrix for the double method.

Finally, we must show that

$$\beta_i = 1 \quad \text{for } i = 0, 1, \dots, \quad (4.2)$$

for the J-CG method. We write (1.1) as

$$\left(\begin{array}{c|c} M_1 & F \\ \hline F^t & M_2 \end{array} \right) \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}. \quad (5.7)$$

Following Concus et al. [2], we choose the initial vector

$$\tilde{x}^{(0)} = \begin{cases} \tilde{x}_1^{(0)} & \text{is some approximation to } \tilde{x}_1 \\ \tilde{x}_2^{(0)} = M_2^{-1}(b_2 - F^t \tilde{x}_1^{(0)}) & , \end{cases} \quad (5.8)$$

so that

$$\tilde{z}_2^{(0)} = 0.$$

As

$$\beta_i = \frac{(\tilde{z}^{(i)}, C\tilde{z}^{(i)})}{(\tilde{z}^{(i)}, A\tilde{z}^{(i)})} \quad \text{for } i = 0, 1, \dots, \quad (\text{Algorithm 3.2})$$

we see that

$$\beta_0 = 1.$$

From the recurrence relation

$$\begin{aligned} \tilde{x}^{(i+1)} &= \alpha_{i+1} [\beta_i (G\tilde{x}^{(i)} + \tilde{k}) + (1 - \beta_i) \tilde{x}^{(i)}] + (1 - \alpha_{i+1}) \tilde{x}^{(i-1)}, \\ &\quad \text{for } i = 0, 1, \dots, \end{aligned} \quad (1.9)$$

where $\alpha_1 = 1$, we get (cf. recurrence relation (2.5))

$$\tilde{z}^{(i+1)} = \alpha_{i+1} [\beta_i (I - C^{-1}A) \tilde{z}^{(i)} + (1 - \beta_i) \tilde{z}^{(i)}] + (1 - \alpha_{i+1}) \tilde{z}^{(i-1)},$$

or equivalently,

$$\begin{aligned} \begin{pmatrix} \tilde{z}_1^{(i+1)} \\ \tilde{z}_2^{(i+1)} \end{pmatrix} &= \alpha_{i+1} \left[\beta_i \begin{pmatrix} 0 & -M_1^{-1}F \\ -M_2^{-1}F^t & 0 \end{pmatrix} \begin{pmatrix} \tilde{z}_1^{(i)} \\ \tilde{z}_2^{(i)} \end{pmatrix} + (1 - \beta_i) \begin{pmatrix} \tilde{z}_1^{(i)} \\ \tilde{z}_2^{(i)} \end{pmatrix} \right] \\ &\quad + (1 - \alpha_{i+1}) \begin{pmatrix} \tilde{z}_1^{(i-1)} \\ \tilde{z}_2^{(i-1)} \end{pmatrix} \quad \text{for } i = 0, 1, \dots \end{aligned} \quad (5.9)$$

It follows that

$$\tilde{z}_1^{(1)} = 0.$$

Thus,

$$\beta_1 = 1.$$

Using (5.9), we readily prove by induction that

$$\begin{cases} \tilde{z}_2^{(2i)} = 0 & \text{and } \beta_{2i} = 1, \\ \tilde{z}_1^{(2i+1)} = 0 & \text{and } \beta_{2i+1} = 1, \end{cases} \quad \text{for } i = 0, 1, \dots \quad (5.10)$$

THEOREM 5.1. The J-CG method with "symmetrization" matrix S defined by (5.6) and the DJ-CG method with "symmetrization" matrix $S(I + G)^{-1}$ are "virtually-equivalent" if the initial vector is chosen by (5.8).

2.6. A Compressed Method and Cyclic Reduction

In this section, we first give two special methods for solving the system

$$\left(\begin{array}{c|c} M_1 & F \\ \hline F^t & M_2 \end{array} \right) \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix} = \begin{pmatrix} \tilde{b}_1 \\ \tilde{b}_2 \end{pmatrix}. \quad (5.7)$$

and then use the result in the last section to show that the two methods are equivalent.

The first method is based on a procedure introduced by Reid [9]. We consider the J-CG method with the special initial vector

$$\tilde{x}^{(0)} = \begin{cases} \tilde{x}_1^{(0)} \\ M_2^{-1}(\tilde{b}_2 - F^t \tilde{x}_1^{(0)}) \end{cases}. \quad (5.8)$$

We have seen that

$$\beta_i = 1 \quad \text{for } i = 0, 1, \dots \quad (4.2)$$

Thus,

$$\begin{aligned} \tilde{x}^{(2i+2)} = & \alpha_{2i+2} \alpha_{2i+1} (G^2 \tilde{x}^{(2i)} + G\tilde{k} + \tilde{k}) + \left[1 + \frac{\alpha_{2i+2}}{\alpha_{2i}} (1 - \alpha_{2i+1}) - \alpha_{2i+2} \right] \tilde{x}^{(2i)} \\ & - \frac{\alpha_{2i+2}}{\alpha_{2i}} (1 - \alpha_{2i+1}) (1 - \alpha_{2i}) \tilde{x}^{(2i-2)} \end{aligned}$$

for $i = 0, 1, \dots$. (4.4)

Since

$$G^2 \tilde{x}^{(2i)} + G\tilde{k} + \tilde{k} - \tilde{x}^{(2i)} = (I + G) \tilde{z}^{(2i)} ,$$

we can simplify (4.4) to

$$\begin{aligned} \tilde{x}^{(2i+2)} - \tilde{x}^{(2i)} \\ = \alpha_{2i+2} \alpha_{2i+1} \left[(I + G) \tilde{z}^{(2i)} + \left(1 - \frac{1}{\alpha_{2i+1}} \right) \left(1 - \frac{1}{\alpha_{2i}} \right) (\tilde{x}^{(2i)} - \tilde{x}^{(2i-2)}) \right] , \end{aligned}$$

for $i = 0, 1, \dots$. (6.1)

To save work, we update only $\tilde{x}_1^{(2i)}$ in (6.1). If we desire $\tilde{x}_2^{(2i)}$, we solve the system

$$M_{22} \tilde{x}_2^{(2i)} = \tilde{b}_2 - F^t \tilde{x}_1^{(2i)} ,$$

as $\tilde{z}_1^{(2i)} = 0$. We update the pseudo-residual vector $\tilde{z}^{(i)}$ recursively by the formula of (5.9)

$$\begin{cases} \tilde{z}_2^{(2i+1)} = \alpha_{2i+1} M_{22}^{-1} F^t \tilde{x}_1^{(2i)} + (1 - \alpha_{2i+1}) \tilde{z}_2^{(2i-1)} , \\ \tilde{z}_1^{(2i+2)} = \alpha_{2i+2} M_{11}^{-1} F \tilde{z}_2^{(2i+1)} + (1 - \alpha_{2i+2}) \tilde{z}_1^{(2i)} , \end{cases} \quad \text{for } i = 0, 1, \dots ,$$

where $\alpha_1 = 1$. We can now present what we call the compressed J-CG method (CJ-CG method).

ALGORITHM 6.1. (CJ-CG method)

(1) Let $\tilde{x}_1^{(0)}$ be an initial approximation to \tilde{x}_1 . Solve the system

$$M_{22} \tilde{x}_2^{(0)} = \tilde{b}_2 - F^t \tilde{x}_1^{(0)} .$$

Compute

$$\hat{z}_1^{(0)} := b_1 - M_1 x_1^{(0)} - F x_2^{(0)}$$

and solve the system

$$M_1 \hat{z}_1^{(0)} = \hat{z}_1^{(0)}.$$

Let

$$i := -1$$

and let $x_1^{(-2)}$ be a vector not equal to $x_1^{(0)}$.

(2) Repeat until $(x_1^{(2i+2)} - x_1^{(2i)}) = 0$:

(a) Let

$$i := i + 1.$$

(b) Compute

$$\alpha_{2i+1} := \begin{cases} 1, & i = 0, \\ \left[1 - \frac{(z_1^{(2i)}, \hat{z}_1^{(2i)})}{(z_2^{(2i-1)}, \hat{z}_2^{(2i-1)})} \right]^{-1}, & i \geq 1, \end{cases}$$

and

$$\hat{z}_2^{(2i+1)} := \alpha_{2i+1} F^t z_1^{(2i)} + (1 - \alpha_{2i+1}) \hat{z}_2^{(2i-1)}.$$

Solve the system

$$M_2 \hat{z}_2^{(2i+1)} = \hat{z}_2^{(2i+1)}.$$

(c) Compute

$$\alpha_{2i+2} := \left[1 - \frac{(z_2^{(2i+1)}, \hat{z}_2^{(2i+1)})}{(z_1^{(2i)}, \hat{z}_1^{(2i)})} \frac{1}{\alpha_{2i+1}} \right]^{-1}$$

and

$$\hat{z}_1^{(2i+2)} := \alpha_{2i+2} F z_2^{(2i+1)} + (1 - \alpha_{2i+2}) \hat{z}_1^{(2i)}.$$

Solve the system

$$M_1 \hat{z}_1^{(2i+2)} = \hat{z}_1^{(2i+2)}.$$

(d) Compute

$$\begin{aligned} & (\tilde{x}_1^{(2i+2)} - \tilde{x}_1^{(2i)}) \\ & := \alpha_{2i+2} \alpha_{2i+1} \left[\tilde{z}_1^{(2i)} + \left(1 - \frac{1}{\alpha_{2i+1}} \right) \left(1 - \frac{1}{\alpha_{2i}} \right) (\tilde{x}_1^{(2i)} - \tilde{x}_1^{(2i-2)}) \right] \end{aligned}$$

and

$$\tilde{x}_1^{(2i+2)} := \tilde{x}_1^{(2i)} + (\tilde{x}_1^{(2i+2)} - \tilde{x}_1^{(2i)}) .$$

(3) Let $\tilde{x}_1^{(2k)}$ be the accepted approximation to \tilde{x}_1 . Solve the system

$$M_{22} \tilde{x}_2^{(2k)} = \tilde{b}_2 - F^t \tilde{x}_1^{(2k)} . \quad \square$$

Let us turn our attention to the double method. It can be viewed as the single method for the matrix equation

$$(I - G^2) \tilde{x} = G \tilde{k} + \tilde{k} . \quad (4.1)$$

As

$$I - G^2 = \left(\begin{array}{c|c} I - M_1^{-1} F M_2^{-1} F^t & 0 \\ \hline 0 & I - M_2^{-1} F^t M_1^{-1} F \end{array} \right) ,$$

the system (4.1) is the uncoupled equations

$$\left\{ \begin{array}{l} (I - M_1^{-1} F M_2^{-1} F^t) \tilde{x}_1 = M_1^{-1} \tilde{b}_1 - M_1^{-1} F M_2^{-1} \tilde{b}_2 \end{array} \right. \quad (6.2a)$$

$$\left\{ \begin{array}{l} (I - M_2^{-1} F^t M_1^{-1} F) \tilde{x}_2 = M_2^{-1} \tilde{b}_2 - M_2^{-1} F^t M_1^{-1} \tilde{b}_1 . \end{array} \right. \quad (6.2b)$$

We refer to the equations of (6.2) as the cyclic reduction (cf. [11]) of the matrix equation

$$\left(\begin{array}{c|c} M_1 & F \\ \hline F^t & M_2 \end{array} \right) \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix} = \begin{pmatrix} \tilde{b}_1 \\ \tilde{b}_2 \end{pmatrix} . \quad (5.7)$$

In the last section, we have chosen the matrix $S(I + G)^{-1}$ as the "symmetrization" matrix for $(I - G^2)$. As

$$[S(I + G)^{-1}] (I - G^2) [S(I + G)^{-1}]^{-1} = S(I - G^2)S^{-1}$$

and

$$S = \left(\begin{array}{c|c} S_1 & 0 \\ \hline 0 & S_2 \end{array} \right) , \quad (5.6)$$

we see that the matrix

$$S_1(I - M_1^{-1}FM_2^{-1}F^t)S_1^{-1}$$

is symmetric and positive definite.

An effective procedure for solving the matrix equation (5.7) consists of

- (a) applying the CG-accelerated single method to solve the reduced system (6.2a) with "symmetrization" matrix S_1 , and
- (b) solving the system

$$M_2 \tilde{x}_2^{(k)} = \tilde{b}_2 - F^t \tilde{x}_1^{(k)} ,$$

where $\tilde{x}_1^{(k)}$ is the accepted approximation to \tilde{x}_1 in (a).

We call this procedure the RS-CG method.

ALGORITHM 6.2 (RS-CG method)

(1) Let

$$x_1^{(0)} := \tilde{x}_1^{(0)} .$$

Solve the system

$$M_1 g^{(0)} = b_1 - FM_2^{-1} b_2 - (M_1 - FM_2^{-1} F^t) x_1^{(0)} .$$

Compute

$$\rho_0 := \frac{(g^{(0)}, M_1 g^{(0)})}{(g^{(0)}, (M_1 - FM_2^{-1} F^t) g^{(0)})}$$

and

$$x_1^{(1)} := \rho_0 g^{(0)} + x_1^{(0)} .$$

Let $r_1 := 1$ and $i := 0$.

(2) Repeat until $g^{(0)} = 0$:

(a) Let

$$i := i + 1 .$$

(b) Solve the system

$$M_1 g_1^{(i)} = b_1 - FM_2^{-1} b_2 - (M_1 - FM_2^{-1} F^t) x_1^{(i)} .$$

(c) Compute

$$\rho_i := \frac{(g^{(i)}, M_1 g^{(i)})}{(g^{(i)}, (M_1 - FM_2^{-1} F^t) g^{(i)})} ,$$

$$r_{i+1} := \left[1 - \frac{\rho_i}{\rho_{i-1}} \frac{(g^{(i)}, M_1 g^{(i)})}{(g^{(i-1)}, M_1 g^{(i-1)})} \frac{1}{r_i} \right]^{-1}$$

and

$$x_1^{(i+1)} := r_{i+1}[\rho_i g^{(i)} + x_1^{(i)}] + (1 - r_{i+1})x_1^{(i-1)}.$$

(3) Let $x_1^{(k)}$ be the accepted approximation to x_1 . Solve the system

$$M_2 x_2^{(k)} = b_2 - F^t x_1^{(k)}. \quad \square$$

We proceed to show how we can derive the RS-CG method from the DJ-CG method. Theorem 5.1 states that the J-CG and DJ-CG methods are "virtually-equivalent," i.e.

$$\tilde{x}^{(2i)} = x^{(i)} \quad \text{for } i = 0, 1, \dots$$

Also we know that

$$\tilde{z}^{(2i)} = 0 \quad \text{for } i = 0, 1, \dots \quad (5.10)$$

Now, let us use the DJ-CG method to update $x_1^{(i)}$:

$$x_1^{(i+1)} = r_{i+1}(\rho_i w_1^{(i)} + x_1^{(i)}) + (1 - r_{i+1})x_1^{(i-1)} \quad \text{for } i = 0, 1, \dots, \quad (6.3)$$

where $r_1 = 1$. As

$$\begin{aligned} w_1^{(i)} &= (I + G)[\underline{k} - (I - G)x_1^{(i)}] \\ &= (I + G)[\underline{k} - (I - G)\tilde{x}^{(2i)}] \\ &= (I + G)\tilde{z}^{(2i)}, \end{aligned}$$

we get

$$\tilde{w}_1^{(i)} = \tilde{z}_1^{(2i)} \quad \text{for } i = 0, 1, \dots \quad (6.4)$$

Let

$$H \equiv I - M_1^{-1} F M_2^{-1} F^t \quad (6.5)$$

and

$$\tilde{h} \equiv M_1^{-1} b_1 - M_1^{-1} F M_2^{-1} b_2. \quad (6.6)$$

Then

$$\begin{aligned} \tilde{z}_1^{(2i)} &= M_1^{-1} b_1 - x_1^{(2i)} - M_1^{-1} F x_2^{(2i)} \\ &= M_1^{-1} b_1 - x_1^{(2i)} - M_1^{-1} F M_2^{-1} (b_2 - F^t x_1^{(2i)}) \\ &= \tilde{h} - (I - H) x_1^{(2i)}. \end{aligned}$$

We define the residual vector $g^{(i)}$ for the reduced system (6.2a) by

$$g^{(i)} \equiv \tilde{h} - (I - H) x_1^{(i)} \quad \text{for } i = 0, 1, \dots \quad (6.7)$$

Thus,

$$g^{(i)} = \tilde{z}_1^{(2i)} = \tilde{w}_1^{(i)} \quad \text{for } i = 0, 1, \dots, \quad (6.8)$$

so that the recurrence relation (6.3) becomes

$$\begin{aligned} x_1^{(i+1)} &= r_{i+1} (\rho_i g^{(i)} + x_1^{(i)}) + (1 - r_{i+1}) x_1^{(i-1)} \\ &\quad \text{for } i = 0, 1, \dots, \end{aligned} \quad (6.9)$$

where $r_1 = 1$. Now,

$$\begin{aligned} (T_{\tilde{w}}^{(i)}, T_{\tilde{w}}^{(i)}) &= (S_{\tilde{z}}^{(2i)}, S_{\tilde{z}}^{(2i)}) \\ &= (S_1 z_1^{(2i)}, S_1 z_1^{(2i)}) \\ &= (g^{(i)}, M_1 g^{(i)}) \end{aligned}$$

and

$$\begin{aligned} (T\tilde{w}^{(i)}, T(I-G^2)\tilde{w}^{(i)}) &= (S\tilde{z}^{(2i)}, S(I-G^2)\tilde{z}^{(2i)}) \\ &= (g^{(i)}, M_1(I-H)g^{(i)}) . \end{aligned}$$

Therefore, the iteration parameters $\{\rho_i\}$ and $\{\gamma_i\}$ in the DJ-CG method can be given by

$$\rho_i = \frac{(g^{(i)}, M_1 g^{(i)})}{(g^{(i)}, M_1(I-H)g^{(i)})} , \quad i \geq 0 , \quad (6.10)$$

and

$$\gamma_{i+1} = \begin{cases} 1 , & i = 0 , \\ \left[1 - \frac{\rho_i}{\rho_{i-1}} \frac{(g^{(i)}, M_1 g^{(i)})}{(g^{(i-1)}, M_1 g^{(i-1)})} \frac{1}{\gamma_i} \right]^{-1} , & i \geq 1 . \end{cases} \quad (6.11)$$

The iterative scheme (6.9) with its iteration parameters given by (6.10) and (6.11) is a second-order procedure for solving the reduced system (6.2a). Since

$$\begin{aligned} (s_1 g^{(i)}, s_1 g^{(j)}) &= (s_1 \tilde{z}_1^{(2i)}, s_1 \tilde{z}_1^{(2j)}) \\ &= (s\tilde{z}^{(2i)}, s\tilde{z}^{(2j)}) \\ &= 0 \quad \text{for } i \neq j, \end{aligned}$$

we conclude that the procedure is identical to the CG-accelerated single method (Theorem 2.1).

We have therefore derived the CJ-CG method from the J-CG method, and the RS-CG method from the DJ-CG method. As the J-CG and the DJ-CG methods are "virtually-equivalent" by Theorem 5.1, we have proved our desired result:

THEOREM 6.1. The CJ-CG and the RS-CG methods are "virtually-equivalent."

2.7. Computer Implementation

There are two popular versions of the conjugate gradient method. One is the original two-term version due to Hestenes and Stiefel [6], the other is a three-term recurrence version due to Rutishauser [10]. We have given our algorithms in the second version for expository purposes. In [8], Reid compared the two versions and found that the Hestenes and Stiefel version is more efficient in both storage and computational work.

Let us present the Hestenes and Stiefel version of the conjugate gradient method for the system

$$A\tilde{x} = \tilde{b} , \quad (1.1)$$

where the matrix A is symmetric and positive definite.

ALGORITHM 7.1. (CG method)

(1) Let $\tilde{x}^{(0)}$ be an initial approximation to \tilde{x} . Compute

$$\tilde{r}^{(0)} := \tilde{b} - A\tilde{x}^{(0)} .$$

Let

$$\tilde{p}^{(0)} := \tilde{r}^{(0)} .$$

Compute

$$\hat{p}^{(0)} := Ap^{(0)},$$

$$\sigma_0 := \frac{(r^{(0)}, \tilde{r}^{(0)})}{(p^{(0)}, \hat{p}^{(0)})}$$

and

$$\tilde{x}^{(1)} := \tilde{x}^{(0)} + \sigma_0 p^{(0)}.$$

Let

$$i := 0.$$

(2) Repeat until $\tilde{r}^{(i)} = 0$ or $p^{(i)} = 0$:

(a) Let

$$i := i + 1.$$

(b) Compute

$$\tilde{r}^{(i)} := \tilde{r}^{(i-1)} - \sigma_{i-1} \hat{p}^{(i-1)}$$

(c) Compute

$$\tau_i := \frac{(r^{(i)}, \tilde{r}^{(i)})}{(\tilde{r}^{(i-1)}, \tilde{r}^{(i-1)})}$$

and

$$p^{(i)} := r^{(i)} + \tau_i \tilde{r}^{(i-1)}.$$

(d) Compute

$$\hat{p}^{(i)} := Ap^{(i)},$$

$$\sigma_i := \frac{(r^{(i)}, \tilde{r}^{(i)})}{(p^{(i)}, \hat{p}^{(i)})}$$

and

$$\tilde{x}^{(i+1)} := \tilde{x}^{(i)} + \sigma_i p^{(i)}.$$

□

Algorithm 7.1 requires one matrix-vector product $A\tilde{p}$ and $5n$ multiplications per iteration. We have to store four vectors \tilde{x} , \tilde{r} , \tilde{p} and $A\tilde{p}$, in addition to the matrix A . The three-term version of the CG method (Algorithm 2.1) requires n more multiplications per iteration and one more vector of storage.

Proceeding in a similar manner as in Section 2.3, we derive from Algorithm 7.1 the two-term version of the preconditioned CG method for the system

$$C^{-1}A\tilde{x} = C^{-1}\tilde{b} . \quad (1.2)$$

ALGORITHM 7.2 (Preconditioned CG method)

(1) Let $\tilde{x}^{(0)}$ be an initial approximation to \tilde{x} . Solve the system

$$C\tilde{z}^{(0)} = \tilde{b} - A\tilde{x}^{(0)} .$$

Let

$$\tilde{p}^{(0)} := \tilde{z}^{(0)}$$

and solve the system

$$C\tilde{p}^{(0)} = A\tilde{p}^{(0)} .$$

Compute

$$\sigma_0 := \frac{(\tilde{z}^{(0)}, s^t s \tilde{z}^{(0)})}{(\tilde{p}^{(0)}, s^t s \tilde{p}^{(0)})}$$

and

$$\tilde{x}^{(1)} := \tilde{x}^{(0)} + \sigma_0 \tilde{p}^{(0)} .$$

Let

$$i := 0 .$$

(2) Repeat until $\underline{z}^{(i)} = \underline{0}$ or $\underline{p}^{(i)} = \underline{0}$:

(a) Let

$$i := i + 1 .$$

(b) Compute

$$\underline{z}^{(i)} := \underline{z}^{(i-1)} - \sigma_{i-1} \bar{\underline{p}}^{(i-1)}$$

(c) Compute

$$\tau_i := \frac{(\underline{z}^{(i)}, S^t S \underline{z}^{(i)})}{(\underline{z}^{(i-1)}, S^t S \underline{z}^{(i-1)})} .$$

(d) Compute

$$\underline{p}^{(i)} := \underline{z}^{(i)} + \tau_i \underline{p}^{(i-1)}$$

and solve the system

$$C \bar{\underline{p}}^{(i)} = A \underline{p}^{(i)} .$$

(e) Compute

$$\sigma_i := \frac{(\underline{z}^{(i)}, S^t S \underline{z}^{(i)})}{(\underline{p}^{(i)}, S^t S \bar{\underline{p}}^{(i)})}$$

and

$$\underline{x}^{(i+1)} := \underline{x}^{(i)} + \sigma_i \underline{p}^{(i)} . \quad \square$$

An iteration of Algorithm 7.2 requires three matrix-vector products $A \underline{p}$, $S^t S \underline{z}$ and $S^t S \bar{\underline{p}}$, one matrix equation solution with coefficient matrix C , and $5n$ multiplications. Storage is required for the three matrices A , C and $S^t S$, and for the five vectors \underline{x} , \underline{z} , \underline{p} , $\bar{\underline{p}}$ and $S^t S \underline{z}$. We can store the vector $S^t S \bar{\underline{p}}$ in $S^t S \underline{z}$.

In the special instance where $S^t S = C$, Algorithm 7.2 reduces to the generalized CG method of Concus et al. [2]:

ALGORITHM 7.3 (Generalized CG method)

(1) Let $\underline{x}^{(0)}$ be an initial approximation to \underline{x} . Compute

$$\underline{r}^{(0)} := \underline{b} - A\underline{x}^{(0)}$$

and solve the system

$$C\underline{z}^{(0)} = \underline{r}^{(0)} .$$

Let

$$\underline{p}^{(0)} := \underline{z}^{(0)} .$$

Compute

$$\sigma_0 := \frac{(\underline{z}^{(0)}, \underline{r}^{(0)})}{(\underline{p}^{(0)}, A\underline{p}^{(0)})}$$

and

$$\underline{x}^{(1)} := \underline{x}^{(0)} + \sigma_0 \underline{p}^{(0)} .$$

Let

$$i := 0 .$$

(2) Repeat until $\underline{z}^{(i)} = \underline{0}$ or $\underline{p}^{(i)} = \underline{0}$:

(a) Let

$$i := i + 1 .$$

(b) Compute

$$\underline{r}^{(i)} := \underline{r}^{(i-1)} - \sigma_{i-1} A\underline{p}^{(i-1)}$$

and solve the system

$$C\underline{z}^{(i)} = \underline{r}^{(i)} .$$

(c) Compute

$$\tau_i := \frac{(\underline{z}^{(i)}, \underline{r}^{(i)})}{(\underline{z}^{(i-1)}, \underline{r}^{(i-1)})} .$$

(d) Compute

$$\underline{p}^{(i)} := \underline{z}^{(i)} + \tau_i \underline{p}^{(i-1)} .$$

(e) Compute

$$\sigma_i := \frac{(\underline{z}^{(i)}, \underline{r}^{(i)})}{(\underline{p}^{(i)}, A \underline{p}^{(i)})}$$

and

$$\underline{x}^{(i+1)} := \underline{x}^{(i)} + \sigma_i \underline{p}^{(i)} . \quad \square$$

Algorithm 7.3 requires one matrix-vector product $A \underline{p}$, one matrix equation solution with coefficient matrix C and $5n$ multiplications for each iteration. We have to store the two matrices A and C , and the four vectors $\underline{x}, \underline{r}, \underline{p}$ and $A \underline{p}$. We can store the vector \underline{z} in $A \underline{p}$.

We now consider solution procedures for the system

$$\left(\begin{array}{c|c} M_1 & F \\ \hline F^t & M_2 \end{array} \right) \begin{pmatrix} \underline{x}_1 \\ \underline{x}_2 \end{pmatrix} = \begin{pmatrix} \underline{b}_1 \\ \underline{b}_2 \end{pmatrix} . \quad (5.7)$$

The J-CG method is Algorithm 7.3 with the preconditioning matrix

$$C = \left(\begin{array}{c|c} M_1 & 0 \\ \hline 0 & M_2 \end{array} \right) . \quad (5.2)$$

As the RS-CG method is the CG method applied to the reduced system (6.2a), we can convert it to the two-term form without much difficulty:

Algorithm 7.4. (RS-CG method)

(1) Let $x_1^{(0)}$ be an initial approximation to x_1 . Compute

$$\hat{g}^{(0)} := b_1 - FM_2^{-1}b_2 - (M_1 - FM_2^{-1}F^t)x_1^{(0)}.$$

Solve the system

$$M_1 g^{(0)} = \hat{g}^{(0)}.$$

Let

$$p_1^{(0)} := g^{(0)}$$

and

$$\hat{p}_1^{(0)} := \hat{g}^{(0)}.$$

Compute

$$\sigma_0 := \frac{(g^{(0)}, \hat{g}^{(0)})}{(p_1^{(0)}, \hat{p}_1^{(0)} - FM_2^{-1}F^t p_1^{(0)})}$$

and

$$x_1^{(1)} := x_1^{(0)} + \sigma_0 p_1^{(0)}.$$

Let

$$i := 0.$$

(2) Repeat until either $p_1^{(i)} = 0$ or $q^{(i)} = 0$:

(a) Let

$$i := i + 1.$$

(b) Compute

$$\hat{g}^{(i)} := \hat{g}^{(i-1)} - \sigma_{i-1} \hat{p}_1^{(i-1)}.$$

(c) Solve the system

$$M_1 q^{(i)} = \hat{g}^{(i)}.$$

(d) Compute

$$\tau_i := \frac{(q^{(i)}, \hat{g}^{(i)})}{(q^{(i-1)}, \hat{g}^{(i-1)})},$$

$$p_1^{(i)} := q^{(i)} + \tau_i p_1^{(i-1)}$$

and

$$\hat{p}_1^{(i)} := \hat{g}^{(i)} + \tau_i \hat{p}_1^{(i-1)}.$$

(e) Compute

$$\sigma_i := \frac{(q^{(i)}, \hat{q}^{(i)})}{(p_1^{(i)}, \hat{p}_1^{(i)} - FM_2^{-1} F^t p_1^{(i)})}$$

and

$$x_1^{(i+1)} := x_1^{(i)} + \sigma_i p_1^{(i)}.$$

(3) Let $x_1^{(k)}$ be the accepted approximation to x_1 . Solve the system

$$M_2 x_2^{(k)} = b_2 - F^t x_1^{(k)}. \quad \square$$

An iteration of Algorithm 7.4 requires one matrix-vector product with F and another product with F^t , one matrix equation solution with coefficient matrix M_1 and another solution with M_2 , and $6n_1$ multiplications, where n_1 is the order of the vector x_1 . We must store the matrix A , the four n_1 -vectors y_1, \hat{g}, p_1 and \hat{p}_1 , and the two $(n-n_1)$ -vectors necessary for computing $FM_2^{-1}F^tp_1$. The vector q can be stored with one of the two $(n-n_1)$ -vectors.

It is not straightforward to convert the CJ-CG method to the two-term form. Since

$$z_2^{(21)} = 0, \quad (5.10)$$

and the method is "virtually-equivalent" to the RS-CG method, we conclude that any two-term version of the CJ-CG method would be essentially Algorithm 7.4. Every two steps of the CJ-CG method (Algorithm 6.1) requires the same number of matrix-vector products and matrix equations solutions as the RS-CG method, plus $3n + 2n_1$ multiplications. Storage is required for the two n -vectors z and \hat{z} , and for the two n_1 -vectors x_1 and δx_1 , in addition to the matrix A .

In the following table, we summarize the work and storage requirements of the four CG-based methods for solving the matrix equation (5.7). We give the necessary work for one step of each method. As one step of the RS-CG method is equivalent to two steps of the other three algorithms, we halve the work requirements for that method.

Method	Matrix-vector products				Solve		Mult.	Storage
	F	F ^t	M ₁	M ₂	M ₁ ⁻¹	M ₂ ⁻¹		
CG	1	1	1	1	0	0	5n	4n
J-CG	1	1	1	1	1	1	5n	4n
CJ-CG	0.5	0.5	0	0	0.5	0.5	1.5n + n ₁	3n
RS-CG	0.5	0.5	0	0	0.5	0.5	3n ₁	2n + 2n ₁

It is obvious that the CJ-CG or the RS-CG method is about twice as fast as the J-CG method. In the next section we give three examples showing how a (block) diagonal matrix preconditioning (\equiv J-CG method) can greatly accelerate the convergence of the CG method. As \tilde{x}_1 can be chosen so that $n_1 \leq n/2$, the RS-CG method is always more efficient than the CJ-CG method in that about n_1 fewer multiplications are required per iteration.

2.8. Applications of the RS-CG Method

The conjugate gradient method is already well established as an effective iterative solution procedure for large and sparse matrix equations. We wish to show here that the CG method, coupled with cyclic reduction, can be particularly useful for solving certain classes of matrix equations which arise in elliptic partial differential equations problems.

Our presentation here is very similar to that of Hageman [5]. We consider the second-order self-adjoint partial differential equation

$$-(a(x,y)u_x)_x - (a(x,y)u_y)_y + \sigma(x,y)u = f(x,y), \quad (x,y) \in R, \quad (8.1)$$

defined in a bounded rectangular region R and subject to the boundary conditions

$$\frac{\partial u(x,y)}{\partial n} = 0 \quad \text{or} \quad u(x,y) = \text{constant}, \quad (x,y) \in \partial R. \quad (8.2)$$

We assume that the given functions a and σ are continuous in $\bar{R} \equiv R \cup \partial R$ with

$$a(x,y) > 0 \quad \text{and} \quad (x,y) \geq 0, \quad (x,y) \in \bar{R}. \quad (8.3)$$

A spatial mesh is then imposed on \bar{R} . Discretizing (8.1) and (8.2) with a five-point difference approximation [11, p. 183], we obtain the matrix equation

$$Ax = b, \quad (1.1)$$

where the $n \times n$ matrix A is symmetric and positive definite.

We divide the mesh points into two sets, one set consisting of square (or black) mesh points and the other consisting of circle (or red) mesh points. The division is done by first making the lower left corner point a circle point and then proceeding by making square points of the four nearest neighbors of the circle points and making circle points of the four nearest neighbors of the square points. We index the mesh points by indexing first all the square points consecutively by rows and then all the circle points consecutively by rows. This ordering is called a point red/black ordering. We give an example in Figure 1.

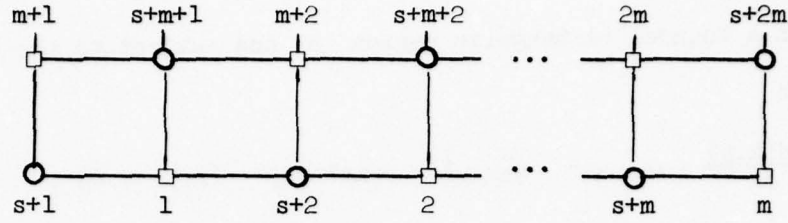


FIGURE 1.

With the point red/black ordering, the matrix equation (1.1) can be written in partitioned form as

$$\left(\begin{array}{c|c} D_1 & H \\ \hline H^t & D_2 \end{array} \right) \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix} = \begin{pmatrix} \tilde{b}_1 \\ \tilde{b}_2 \end{pmatrix}, \quad (8.4)$$

where D_1 and D_2 are diagonal matrices. We have shown in Section 2.6 that the system (8.4) can be cyclically reduced to the lower order matrix equations

$$\begin{cases} (I - D_1^{-1} H D_2^{-1} H^t) \tilde{x}_1 = D_1^{-1} \tilde{b}_1 - D_1^{-1} H D_2^{-1} \tilde{b}_2, & (8.5a) \\ (I - D_2^{-1} H^t D_1^{-1} H) \tilde{x}_2 = D_2^{-1} \tilde{b}_2 - D_2^{-1} H^t D_1^{-1} \tilde{b}_1. & (8.5b) \end{cases}$$

The two sets of equations are disjoint, one set involving only the square points and the other involving only the circle points. Let us consider the system (8.5a). Because of our use of the five-point difference approximation, the matrix $(I - D_1^{-1} H D_2^{-1} H^t)$ corresponds to a nine-point difference approximation. Figure 2 illustrates the typical coupling of a (shaded) square mesh point to the eight (darkened) neighboring square mesh points.

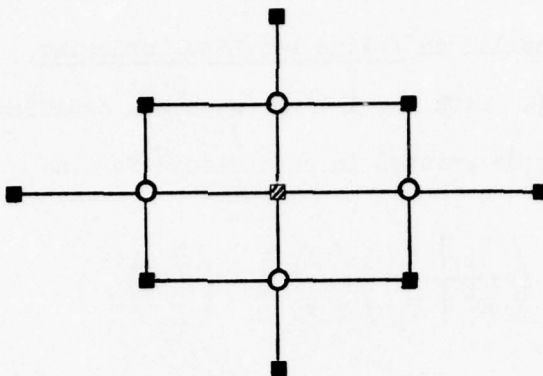


FIGURE 2

We proceed to describe another way to order the mesh points. Assume that there are $p\ell$ horizontal mesh lines. We split the p blocks of successive ℓ horizontal mesh lines into square blocks and circle blocks; we make the first block a circle block, the second one a square block, the third a circle block, and so on. We index the blocks from 1 to p by indexing first all the square blocks consecutively and then all the circle blocks consecutively, as shown in Figure 3.

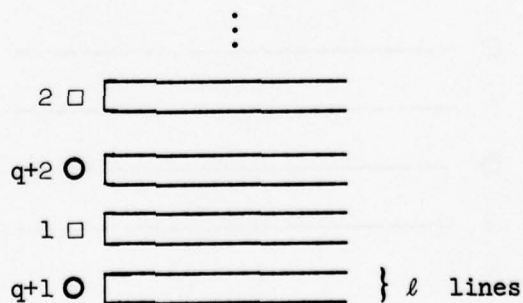


FIGURE 3.

Within each block, the mesh points are indexed consecutively by rows.

This ordering is called an ℓ -line red/black ordering.

Let $\ell = 1$. With the 1-line red/black ordering, the matrix equation (1.1) can be written in partitioned form as

$$\left(\begin{array}{c|c} T_1 & H \\ \hline H^t & T_2 \end{array} \right) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}, \quad (8.6)$$

where T_1 and T_2 are tridiagonal matrices. The cyclically reduced systems obtainable from (8.6) are

$$(I - T_1^{-1} H T_2^{-1} H^t) x_1 = T_1^{-1} b_1 - T_1^{-1} H T_2^{-1} b_2, \quad (8.7a)$$

$$(I - T_2^{-1} H^t T_1^{-1} H) x_2 = T_2^{-1} b_2 - T_2^{-1} H^t T_1^{-1} b_1. \quad (8.7b)$$

The system (8.7a) involves only the square mesh lines. Indeed, the matrix $(I - T_1^{-1} H T_2^{-1} H^t)$ corresponds to a three-line difference approximation, as illustrated in Figure 4.

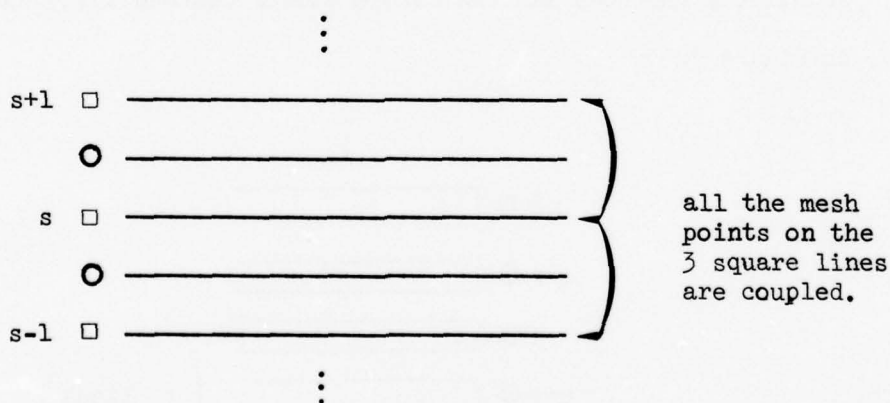
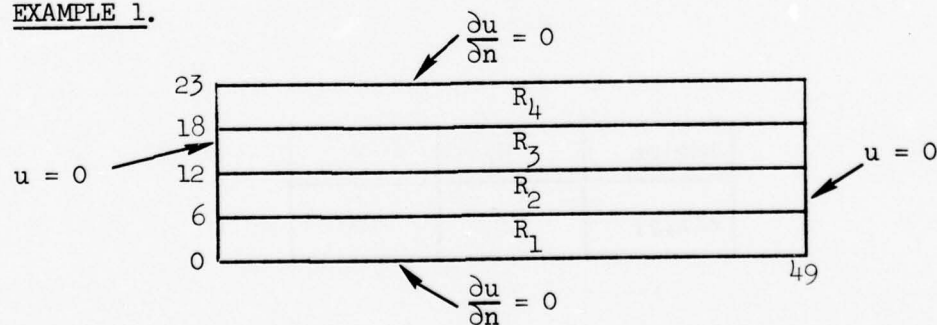


FIGURE 4.

We wish to compare the CG, J-CG and RS-CG methods for the numerical solution of the differential equation (8.1). The CJ-CG method is not included because it is both "virtually-equivalent" to and less efficient than the RS-CG method. We have written three computer programs in the ALGOL W language implementing the three methods. Our programs were run on an IBM 370/168 computer at the Stanford Linear Accelerator Center. The time we give is machine execution time in seconds.

We choose the functions a , σ and f , and the boundary conditions so that the exact solution to the discretized problem is known. Also, we choose the initial vector $\tilde{x}^{(0)}$ so that each component of the vector is a random number from a uniform distribution in the open interval $(-10^3, 10^3)$. Each method has a different initial vector. The iterative procedure terminates when the maximum norm of the error vector is less than 10^{-3} .

EXAMPLE 1.



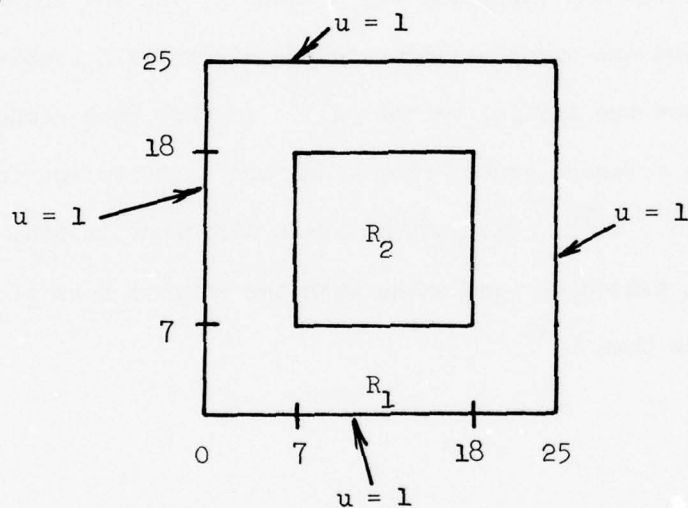
Region	R_1	R_2	R_3	R_4
$a(x,y)$	5	4	9	1

We choose $\sigma(x,y) \equiv 0$ and $f(x,y) \equiv 0$ so that the null vector $\underline{0}$ is the solution to the discretized problem. We use the point red/black ordering.

$$n = 1152$$

Method	Iter	Time
CG	262	42.61
J-CG	136	25.61
RS-CG	66	9.19

EXAMPLE 2.



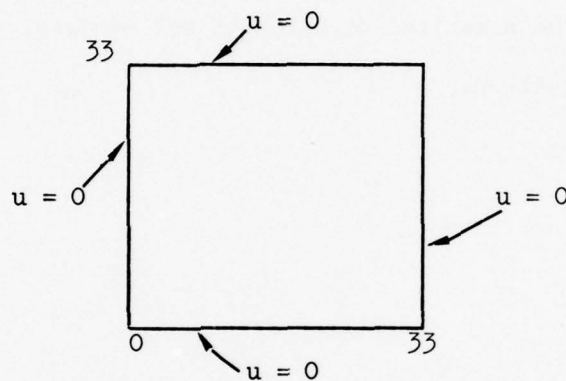
Region	R_1	R_2
$a(x,y)$	2	500

We choose $\sigma(x,y) \equiv 0.05$ and $f(x,y) \equiv 0.05$ so that the vector $\underline{e} \equiv (1,1,\dots,1)^t$ is the solution to the discretized problem. We use the point red/black ordering.

$n = 576$

Method	Iter	Time
CG	400	33.34
J-CG	78	7.52
RS-CG	40	2.80

EXAMPLE 3.



We choose $a(x,y) \equiv 1$, $\sigma(x,y) \equiv 0$ and $f(x,y) \equiv 0$ so that the null vector solves the discretized problem. We use a 1-line red/black ordering. The tridiagonal systems of linear equations are solved by Gaussian elimination; pivoting is not necessary because the systems are positive definite.

$n = 1024$

Method	Iter	Time
CG	86	12.89
J-CG	62	12.05
RS-CG	34	4.55

Our examples show how a (block) diagonal matrix preconditioning can greatly improve the convergence rate of the CG method, especially for the case where the diagonal elements of the coefficient matrix are of different orders of magnitude (see Example 2). The RS-CG method, as expected, requires about half as many iterations as the J-CG method. The saving in machine execution time is, however, even more substantial. It is thus fair to conclude that the RS-CG method is an effective procedure for the numerical solution of self-adjoint elliptic partial differential equations.

REFERENCES

- [1] Axelsson, O., "Solution of linear systems of equations: iterative methods," Chapter 1 of Sparse Matrix Techniques, Copenhagen 1976, Springer-Verlag, Berlin (1977).
- [2] Concus, P., Golub, G.H., and O'Leary, D.P., "A generalized conjugate gradient method for the numerical solution of elliptic partial differential equations," in Sparse Matrix Computations (J. R. Bunch and D. J. Rose, Ed.), Academic Press, New York (1976), 309-332.
- [3] Faddeev, D.K., and Faddeeva, V.N., Computational Methods of Linear Algebra, Freeman, San Francisco (1963).
- [4] Golub, G.H., and Varga, R.S., "Chebyshev semi-iterative methods, successive over-relaxation iterative methods and second order Richardson iterative methods," Numer. Math. 3 (1961), 147-168.
- [5] Hageman, L.A., "Block iterative methods for two-cyclic matrix equations with special application to the numerical solution of the second-order self-adjoint elliptic partial differential equations in two dimensions," Report WAPD-TM-327, Bettis Atomic Power Laboratory (1962).
- [6] Hestenes, M., and Stiefel, E., "Methods of conjugate gradients for solving linear equations," J. Research NBS 49 (1952), 409-436.
- [7] Meijerink, J.A., and van der Vorst, H.A., "An iterative solution method for linear systems of which the coefficient matrix is a symmetric M-matrix," Math. Comp. 31 (1977), 148-162.
- [8] Reid, J.K., "On the method of conjugate gradients for the solution of large sparse systems of linear equations," in Large Sparse Sets of Linear Equations, (J. K. Reid, Ed.), Academic Press, New York (1971), 231-254.
- [9] Reid, J.K., "The use of conjugate gradients for systems of linear equations possessing 'Property A'," SIAM J. Numer. Anal. 9 (1972), 325-332.
- [10] Rutishauser, H., "Theory of gradient methods," Chapter 2 of Refined Iterative Methods for Computation of the Solution and the Eigenvalues of Self-adjoint Boundary Value Problems, M. Engeli, Th. Ginsburg, H. Rutishauser and E. Stiefel, Birkhauser, Basel (1959).
- [11] Varga, R.S., Matrix Iterative Analysis, Prentice-Hall, New York (1962).

III. QUADRATIC PROGRAMMING WITH M-MATRICES

3.1. Introduction

In this chapter, we address the lower and upper bounds quadratic program

$$\begin{aligned} \min_{\tilde{x}} \quad & \frac{1}{2} \tilde{x}^t A \tilde{x} - \tilde{x}^t b \\ \text{subject to} \quad & \underline{c} \leq \tilde{x} \leq \underline{d}, \end{aligned} \tag{1.1}$$

where A is an $n \times n$ M-matrix and \underline{b} , \underline{c} and \underline{d} are n -vectors. An important special case of (1.1) is the linear complementarity problem, in which $\underline{c} = \underline{0}$ and $\underline{d} = \infty$:

$$\begin{aligned} \min_{\tilde{x}} \quad & \frac{1}{2} \tilde{x}^t A \tilde{x} - \tilde{x}^t b \\ \text{subject to} \quad & \tilde{x} \geq \underline{0}. \end{aligned} \tag{1.2}$$

We assume that the matrix A is large and sparse. The problems (1.1) and (1.2) find applications in the numerical solution of free boundary problems for elliptic partial differential equations. Such problems include various types of Dirichlet problems with obstacles ([7] and [10]), and models of the journal bearing [5] and of the application of torsion to a bar [1].

We define an M-matrix as follows.

DEFINITION 1.1 [11, p. 85]. A real square matrix $A = (a_{ij})$ with $a_{ij} \leq 0$ for all $i \neq j$ is an M-matrix if A is nonsingular, and $A^{-1} \geq 0$.

We lose no generality by restricting our attention to symmetric M-matrices, for we can replace the matrix A by its symmetric part $(A + A^t)/2$ in the quadratic form of (1.1) or (1.2) without changing the value of the quadratic form. The next lemma states that a symmetric M-matrix is positive definite. Thus, the problem (1.1) or (1.2) always has a unique solution.

DEFINITION 1.2 [11, p. 85]. A real square matrix $A = (a_{ij})$ with $a_{ij} \leq 0$ for all $i \neq j$ is a Stieltjes matrix if A is symmetric and positive definite.

LEMMA 1.1 [11, p. 85 and 87]. A symmetric M-matrix is a Stieltjes matrix and vice versa.

There are many good algorithms for solving problems (1.1) and (1.2) when the matrix A is positive definite (see [4]). However, it is possible to exploit the special properties of an M-matrix to obtain more efficient special algorithms. Chandrasekaran [2] proposed an algorithm for the linear complementarity problem (1.2), and Pang [8] developed an analogous algorithm for the lower and upper bounds problem (1.1). Cottle and Goheen [4] performed an extensive comparison of the latter method with four other well known algorithms. Their results indicate that Pang's method is the most efficient. They also described a preprocessing scheme that could be used with the method. Their scheme would identify some of the variables which will be at their bounds in

the optimal solution. Those variables could then be eliminated from further consideration.

In this chapter, we study the algorithms of Chandrasekaran and Pang. We demonstrate the special monotone behavior of the iterate and gradient vectors. The result on the gradient vector is new. It leads us to consider a simple updating procedure which preserves the monotonicity of both vectors. The procedure uses the fact that an M-matrix has a nonnegative inverse. Two new algorithms are then constructed by incorporating this updating procedure into the two described algorithms. Numerical tests show that our new algorithms can be twice as fast as the original methods (see Section 3.5).

We also consider the extension of problem (1.1) to the case when some components of the bounds are infinite. We show how we can compute finite a priori bounds on all the variables and reduce the extended problem to one with finite bounds (see Section 3.4). This result is new.

Let

$$v(\underline{x}) \equiv \frac{1}{2} \underline{x}^t A \underline{x} - \underline{x}^t \underline{b} \quad (1.3)$$

and

$$\underline{g}(\underline{x}) \equiv A \underline{x} - \underline{b} , \quad (1.4)$$

i.e., $v(\underline{x})$ is the value of the quadratic form at \underline{x} and $\underline{g}(\underline{x})$ is the gradient vector of the quadratic form at \underline{x} . The solution \underline{x} for problem (1.1) must satisfy the optimality conditions:

$$\begin{array}{ll} \text{if } x_i = c_i & \text{then } g_i(\underline{x}) \geq 0 \\ \text{if } c_i < x_i < d_i & \text{then } g_i(\underline{x}) = 0 \\ \text{if } x_i = d_i & \text{then } g_i(\underline{x}) \leq 0 , \end{array} \quad (1.5)$$

for $i = 1, 2, \dots, n$. The solution \underline{x} for problem (1.2) must satisfy the optimality conditions:

$$\begin{array}{ll} \text{if } x_i = 0 & \text{then } g_i(\underline{x}) \geq 0 \\ \text{if } x_i > 0 & \text{then } g_i(\underline{x}) = 0, \end{array} \quad (1.6)$$

for $i = 1, 2, \dots, n$.

The following lemma and notations are used in this chapter.

DEFINITION 1.3 [11, p. 30]. A principal submatrix of an $n \times n$ matrix A is any matrix obtained by crossing out any j rows and the corresponding j columns of A , where $1 \leq j < n$.

LEMMA 1.2 [11, p. 86]. Any principal submatrix of an M -matrix is an M -matrix.

NOTATIONS.

- (1) $N \equiv \{1, 2, \dots, n\}$.
- (2) $\bar{J} \equiv N \setminus J$ for any subset J of N , i.e. \bar{J} is the complement of J in N .
- (3) $\underline{x}_J \equiv (x_{j_1}, x_{j_2}, \dots, x_{j_p})$, where
 $J = \{j_1, j_2, \dots, j_p\} \subseteq N$ and $j_1 < j_2 < \dots < j_p$.
- (4) $A_{JK} \equiv (c_{rs})$, where
 $c_{rs} = a_{j_r k_s}$, $J = \{j_1, j_2, \dots, j_p\} \subseteq N$,
 $K = \{k_1, k_2, \dots, k_q\} \subseteq N$, $j_1 < j_2 < \dots < j_p$ and $k_1 < k_2 < \dots < k_q$.

3.2. Linear Complementarity Problem

In this section, we consider the linear complementarity problem

$$\begin{aligned} \min_{\tilde{x}} \quad & \frac{1}{2} \tilde{x}^t A \tilde{x} - \tilde{x}^t \tilde{b} \\ \text{subject to} \quad & \tilde{x} \geq \underline{0}, \end{aligned} \tag{1.2}$$

where A is a symmetric M-matrix. We assume that \tilde{b} contains both positive and negative components to eliminate two trivial cases.

LEMMA 2.1.

- (1) If $\tilde{b} \leq \underline{0}$, then $\tilde{x} = \underline{0}$ solves (1.2).
- (2) If $\tilde{b} \geq \underline{0}$, then $\tilde{x} = A^{-1} \tilde{b}$ solves (1.2).

Proof.

- (1) Assume \tilde{x} solves (1.2) with $\tilde{x}_P > \underline{0}$ and $\tilde{x}_{\bar{P}} = \underline{0}$ for some nonempty set $P \subseteq N$. The optimality conditions (1.6) give

$$g_P(\tilde{x}) \equiv A_{PP} \tilde{x}_P - \tilde{b}_P = \underline{0}.$$

Hence

$$\tilde{x}_P = A_{PP}^{-1} \tilde{b}_P \leq \underline{0} \quad \text{because} \quad A_{PP}^{-1} \geq 0 \quad \text{and} \quad \tilde{b}_P \leq \underline{0}.$$

Contradiction.

- (2) We have

$$\tilde{x} = A^{-1} \tilde{b} \geq \underline{0} \quad \text{because} \quad A^{-1} \geq 0 \quad \text{and} \quad \tilde{b} \geq \underline{0},$$

and

$$g(\tilde{x}) = A \tilde{x} - \tilde{b} = \underline{0}.$$

Thus, \tilde{x} solves (1.2). \square

We proceed to prove two more lemmas and an important theorem.

LEMMA 2.2. Suppose we have a Stieltjes matrix A and a nonempty subset I of N . The matrix

$$V \equiv A_{II} - A_{II} A_{II}^{-1} A_{II}$$

is a Stieltjes matrix.

Proof. The matrix V is symmetric because $A_{II} = A_{II}^t$. Let $W \equiv A_{II}^{-1}$. Since A is a positive definite matrix, so is W . A fortiori, so are W_{II} and W_{II}^{-1} . But $V = W_{II}^{-1}$ [6, p. 99]. The matrix A_{II} is an M-matrix by Lemma 1.2. Hence $A_{II}^{-1} \geq 0$. The off-diagonal elements of V are therefore nonpositive because the off-diagonal elements of A_{II} are nonpositive, and $A_{II} \leq 0$, $A_{II}^{-1} \geq 0$ and $A_{II} \leq 0$. \square

LEMMA 2.3. Suppose we have vectors \tilde{x} and \tilde{y} with $\tilde{x}_J = \tilde{y}_J$ and $\tilde{x}_{\bar{J}} \leq \tilde{y}_{\bar{J}}$ for some nonempty subset J of N . Then $g_J(\tilde{x}) \geq g_J(\tilde{y})$.

Proof.

$$\begin{aligned} g_J(\tilde{x}) &= A_{JJ} \tilde{x}_J + A_{J\bar{J}} \tilde{x}_{\bar{J}} - b_J \\ &= A_{JJ} \tilde{y}_J + A_{J\bar{J}} \tilde{y}_{\bar{J}} - b_J + A_{J\bar{J}} (\tilde{x}_{\bar{J}} - \tilde{y}_{\bar{J}}) \\ &= g_J(\tilde{y}) + A_{J\bar{J}} (\tilde{x}_{\bar{J}} - \tilde{y}_{\bar{J}}) \\ &\geq g_J(\tilde{y}) \end{aligned}$$

because $A_{J\bar{J}} \leq 0$ and $\tilde{x}_{\bar{J}} - \tilde{y}_{\bar{J}} \leq 0$. \square

THEOREM 2.1. Assume $A_{JJ}^{-1}b_J > 0$ for some nonempty subset J of N .

Define a vector \tilde{x} with $\tilde{x}_J = A_{JJ}^{-1}b_J$ and $\tilde{x}_{\bar{J}} = 0$. Let K and L be two sets partitioning \bar{J} so that $g_K(\tilde{x}) < 0$ and $g_L(\tilde{x}) \geq 0$.

If the set K is empty, then

(1) \tilde{x} solves (1.2),

else

(2) let $Q := J \cup K$. Construct \tilde{y} with $\tilde{y}_Q = A_{QQ}^{-1}b_Q$ and $\tilde{y}_L = 0$.

We get

(a) $\tilde{y}_J \geq \tilde{x}_J > 0$, $\tilde{y}_K > 0$,

(b) $g_L(\tilde{y}) \leq g_L(\tilde{x})$,

and

(c) $v(\tilde{y}) < v(\tilde{x})$. (Remember $v(\tilde{z}) \equiv \frac{1}{2} \tilde{z}^t A \tilde{z} - \tilde{z}^t b$).

Proof.

(1) The optimality conditions (1.6) are satisfied because $g_J(\tilde{x}) = 0$ and $g_{\bar{J}}(\tilde{x}) \geq 0$.

(2a) From the system

$$A_{JJ}\tilde{y}_J + A_{JK}\tilde{y}_K = b_J,$$

$$A_{KJ}\tilde{y}_J + A_{KK}\tilde{y}_K = b_K,$$

we obtain

$$\tilde{y}_K = (A_{KK} - A_{KJ}A_{JJ}^{-1}A_{JK})^{-1} (b_K - A_{KJ}A_{JJ}^{-1}b_J).$$

But

$$\begin{aligned} b_K - A_{KJ}A_{JJ}^{-1}b_J &= b_K - A_{KJ}\tilde{x}_J \\ &= -g_K(\tilde{x}) > 0 \end{aligned}$$

and

$$(A_{KK} - A_{KJ} A_{JJ}^{-1} A_{JK})^{-1} \geq 0 \quad \text{by Lemma 2.2.}$$

Thus,

$$x_K > 0.$$

Also,

$$\begin{aligned} x_J &= A_{JJ}^{-1} b_J - A_{JJ}^{-1} A_{JK} x_K \\ &\geq A_{JJ}^{-1} b_J \quad \text{because } A_{JJ}^{-1} \geq 0, A_{JK} \leq 0 \text{ and } x_K > 0 \\ &= \bar{x}_J \\ &> 0. \end{aligned}$$

(2b) We have shown in (2a) that $x_J \geq \bar{x}_J$ and $x_K > 0 = \bar{x}_K$. Let $\bar{L} = J \cup K$ and $x_{\bar{L}} = \bar{x}_{\bar{L}} = 0$. The conditions of Lemma 2.3 are satisfied and therefore $g_L(x) \leq g_L(\bar{x})$.

$$\begin{aligned} (2c) \quad 2 \cdot v(x) &= x_Q^t A_{QQ} x_Q - 2 x_Q^t b_Q \\ &= -b_Q^t A_{QQ}^{-1} b_Q + (x_Q - A_{QQ}^{-1} b_Q)^t A_{QQ} (x_Q - A_{QQ}^{-1} b_Q) - b_Q^t A_{QQ}^{-1} b_Q \\ &= \bar{x}_Q^t A_{QQ} \bar{x}_Q - 2 \bar{x}_Q^t b_Q \\ &= 2 \cdot v(\bar{x}). \quad \square \end{aligned}$$

Let Z denote the index set of the constrained (zero) components of the iterate vector and P the index set of the unconstrained (positive) components. We can describe the lower bound algorithm in [2] as follows.

ALGORITHM 2.1 (Chandrasekaran's method). Let

$$P := \{i \in N \mid b_i > 0\} \quad \text{and} \quad Z := \bar{P}.$$

Define \tilde{x} so that $\tilde{x}_P = A_{PP}^{-1}b_P$ and $\tilde{x}_Z = 0$. Let J be a nonempty set.

Repeat until either set Z or J is empty:

1. Compute $g_Z := A_{ZP}\tilde{x}_P - b_Z$
2. Let $J := \{j \in Z \mid g_j < 0\}$.
3. If set J is nonempty, then
 - (a) Let $P := P \cup J$ and $Z := Z \setminus J$.
 - (b) Reconstruct \tilde{x} so that $\tilde{x}_P = A_{PP}^{-1}b_P$ and $\tilde{x}_Z = 0$. \square

Part (2a) of Theorem 2.1 says that the iterate vector is non-decreasing in value. Thus, a positive component stays positive and no element leaves the index set P . We enlarge set P when we release variables with negative gradients from their constraints. If no such variables exist, then we have computed the solution (Part (1) of Theorem 2.1). Hence Algorithm 2.1 always terminates. It is a descent method by Part (2c) of the same theorem.

Most of the work in an iteration of Algorithm 2.1 is spent in the solution of a matrix equation. We can use an iterative method. The conjugate gradient method (cf. [9]) or its generalized variant (cf. [3]) is particularly effective for solving a large sparse set of linear equations when the coefficient matrix is symmetric and positive definite. Furthermore, we can construct a good starting vector from the solution to the matrix equation in the previous iteration. A straightforward strategy is to place zeros in the new positions. In many cases, the quadratic program (1.2) arises from an approximation to a continuous problem and some sort of interpolation and extrapolation schemes can be successfully exploited.

It is obvious that we do not need to solve the matrix equations to full accuracy in any but the last iteration. This gives another reason for using an iterative method instead of a direct method for solving the matrix equations.

The intent of solving the matrix equation is to go from the minimum of one subspace S_1 to the minimum of a larger subspace S_2 that contains S_1 . The computation of the minima increases the size of the index set P rapidly. A disadvantage is that the cost involved can be substantial. Thus, we are interested in the possibility of using some simple computing process that lets us enlarge the set P with very little work. It will take more iterations to determine the final index set P , but the total cost may be lower because of the smaller amount of work per iteration.

Suppose we are at the start of an iteration of Algorithm 2.1. Let \underline{x} and $\underline{g}(\underline{x})$ be the iterate and gradient vectors, respectively. Let P be the index set of positive components of \underline{x} and J be the index set of zero components with negative gradients, i.e.,

$$P := \{j \in \mathbb{Z}_n \mid x_j > 0\} \quad \text{and} \quad J := \{j \in \bar{P} \mid g_j < 0\}.$$

Assume set J is nonempty. Let $Q := P \cup J$. We construct the new iterate $\hat{\underline{x}}$ so that $\hat{\underline{x}}_{\bar{Q}} = \underline{x}_{\bar{Q}}$ and $\hat{\underline{x}}_Q = \underline{y}$, where \underline{y} solves the matrix equation

$$A_{QQ}\underline{y} = \underline{b}_Q. \quad (2.1)$$

Let us consider a cheaper way to construct $\hat{\underline{x}}$: update only the j -th components, where $j \in J$. In other words, $\hat{\underline{x}}_{\bar{J}} = \underline{x}_{\bar{J}}$ and $\hat{\underline{x}}_J = \underline{z}$, where \underline{z} solves the matrix equation

$$B\underline{z} = -\underline{g}_J \quad (2.2)$$

and B is an M -matrix of order equal to the size of J . Thus $\underline{z} > \underline{0}$. We can choose B as A_{JJ} , but a better strategy is to choose B so that the matrix equation (2.2) is easily solvable. Regardless of which B we choose, it is very important that the new iterate $\hat{\underline{x}}$ satisfies the inequality

$$\hat{\underline{x}}_{\bar{Q}} \leq A_{QQ}^{-1} \underline{b}_Q. \quad (2.3)$$

If (2.3) does not hold, then the gradients of some constrained variables

may have decreased too much in value (cf. Lemma 2.3) and those variables may be erroneously released from their constraints.

We propose to choose B as C_{JJ} , where C is an $n \times n$ M-matrix so that $C \geq A$. The next theorem shows that the inequality (2.3) always holds for these choices of B and C . We now give a simple way to construct the matrix C .

LEMMA 2.4 [11, p. 85]. Let A be an M-matrix, and let C be any matrix obtained from A by setting certain off-diagonal entries of the matrix to zero. Then, C is also an M-matrix.

A possible choice for C is therefore a (block) diagonal matrix with the same (block) diagonal part as A .

ASSUMPTION 2.1. We have chosen an M-matrix C such that $C \geq A$ and that matrix equations with C or a principal submatrix of C as the coefficient matrix are easily solvable.

LEMMA 2.5. Suppose there is a vector $\tilde{x} \geq 0$ such that $\tilde{x}_J = 0$ and $g_J(\tilde{x}) < 0$ for some nonempty subset J of N . Construct a vector y so that $y_{\bar{J}} = \tilde{x}_{\bar{J}}$ and $y_J = -C_{JJ}^{-1}g_J(\tilde{x})$. Then $g_J(y) \leq 0$.

Proof. C_{JJ} is an M-matrix by Lemma 1.2. Hence $C_{JJ}^{-1} \geq 0$. But $A_{JJ} \leq C_{JJ}$. Thus,

$$A_{JJ}C_{JJ}^{-1} \leq C_{JJ}C_{JJ}^{-1} = I,$$

which implies

$$I - A_{JJ}C_{JJ}^{-1} \geq 0.$$

Finally,

$$\begin{aligned}
g_J(y) &= A_{JJ}y_J + A_{J\bar{J}}y_{\bar{J}} - b_J \\
&= -A_{JJ}C_{JJ}^{-1}g_J(x) + A_{J\bar{J}}x_{\bar{J}} - b_J \\
&= -A_{JJ}C_{JJ}^{-1}g_J(x) + g_J(x) \quad \text{because} \quad x_J = 0 \\
&= (I - A_{JJ}C_{JJ}^{-1}) g_J(x) \\
&\leq 0 \quad \text{because} \quad g_J(x) < 0. \quad \square
\end{aligned}$$

LEMMA 2.6. Suppose there is a nonempty subset J of N and two vectors x and y such that $x_{\bar{J}} = 0$, $g_J(x) = 0$, $y_{\bar{J}} = 0$ and $g_J(y) \leq 0$. Then $x_J \geq y_J$.

Proof.

$$\begin{aligned}
A_{JJ}(x_J - y_J) &= A_{JJ}x_J - b_J - (A_{JJ}y_J - b_J) \\
&= g_J(x) - g_J(y) \\
&\geq 0.
\end{aligned}$$

Hence

$$x_J - y_J \geq 0. \quad \square$$

THEOREM 2.2. Assume we have a vector $x \geq 0$. Let $P := \{j \in N \mid x_j > 0\}$ and $J := \{j \in \bar{P} \mid g_j(x) < 0\}$. Suppose that both index sets P and J are nonempty, and that $g_P(x) \leq 0$. Construct a vector y such that $y_J = -C_{JJ}^{-1}g_J(x)$ and $y_{\bar{J}} = x_{\bar{J}}$. Let $K := \bar{P} \setminus J$. Then

- (1) $g_P(y) \leq g_P(x) \leq 0$, $g_J(y) \leq 0$ and $g_K(y) \leq g_K(x)$.
- (2) $x_Q \leq A_{QQ}^{-1} b_Q$, where $Q = P \cup J$.
- (3) $v(y) < v(x)$.

Proof.

- (1) By Lemma 2.5, $g_J(y) \leq 0$. Since $x_J > 0 = x_{\bar{J}}$ and $y_{\bar{J}} = x_{\bar{J}}$, we get $g_{\bar{J}}(y) \leq g_{\bar{J}}(x)$ by Lemma 2.3. But $\bar{J} = P \cup K$.
- (2) Define a vector z so that $z_Q = A_{QQ}^{-1} b_Q$ and $z_{\bar{Q}} = y_{\bar{Q}} = 0$. Hence $g_Q(z) = 0$. From Part (1), $g_Q(y) \leq 0$. Thus, $z_Q \geq x_Q$ by Lemma 2.6.
- (3)
$$\begin{aligned} v(y) - v(x) &= \frac{1}{2} x_J^t A_{JJ} x_J + x_J^t A_{J\bar{J}} x_{\bar{J}} - x_J^t b_J \\ &= \frac{1}{2} x_J^t A_{JJ} (-C_{JJ}^{-1} g_J(x)) + x_J^t (A_{J\bar{J}} x_{\bar{J}} - b_J) \\ &= x_J^t (I - \frac{1}{2} A_{JJ} C_{JJ}^{-1}) g_J(x). \end{aligned}$$

We have shown that $I - A_{JJ} C_{JJ}^{-1} \geq 0$. Hence $(I - (A_{JJ} C_{JJ}^{-1})/2) \geq 0$ and equality is not possible. Since $x_J > 0$ and $g_J(x) < 0$, we get $v(y) < v(x)$. \square

We now have the tool to modify Algorithm 2.1. Instead of going from one constrained minimum to another through solving matrix equations involving the unconstrained variables, we take descent steps through solving very simple matrix equations involving only those just released variables. We call our technique "partial updating". It lowers the

gradients of some constrained variables and we can again release those variables with negative gradients from the lower bound using the same technique. The process is repeated until the gradients of the constrained variables are all nonnegative. Then we compute the constrained minimum by solving the matrix equation

$$A_{PP}x = b_P \quad (2.4)$$

and defining the new iterate \tilde{x} so that $\tilde{x}_P = x$ and $\tilde{x}_{\bar{P}} = \tilde{Q}$. The iterate vector from the "partial updating" is usually a very good initial vector for solving (2.4) using an iterative method.

ALGORITHM 2.2. (Modified Chandrasekaran's method). Let

$$J := \{j \in N | b_j > 0\}, \quad P := J,$$

$$Z := \bar{J}, \quad \tilde{x} = \tilde{Q} \quad \text{and} \quad g := -b.$$

Repeat until either set Z or J is empty:

1. Repeat until either set Z or J is empty:

$$(a) \quad \text{Compute } \tilde{x}_J := -C_{JJ}^{-1}g_J.$$

$$(b) \quad \text{Update } g_Z := g_Z + A_{ZJ}\tilde{x}_J.$$

$$(c) \quad \text{Redefine } J := \{j \in Z | g_j < 0\}.$$

$$(d) \quad \text{Let } P := P \cup J \quad \text{and} \quad Z := Z \setminus J.$$

$$2. \quad \text{Compute } \tilde{x}_P := A_{PP}^{-1}b_P.$$

$$3. \quad \text{Compute } g_Z := A_{ZP}\tilde{x}_P - b_Z.$$

$$4. \quad \text{Redefine } J := \{j \in Z | g_j < 0\}. \quad \text{Let } P := P \cup J \quad \text{and} \quad Z := Z \setminus J. \quad \square$$

We can easily extend Algorithm 2.2 to solve the general lower-bound quadratic program

$$\min_{\tilde{x}} \frac{1}{2} \tilde{x}^t A \tilde{x} - \tilde{x}^t \tilde{b} \quad (2.5)$$

subject to $\tilde{x} \geq \tilde{c}$.

ALGORITHM 2.3. Let

$$P := \emptyset, \quad Z := N, \quad \tilde{x} = \tilde{c}, \quad \text{and} \quad \tilde{g} := A\tilde{c} - \tilde{b}.$$

Set $J := \{j \in Z \mid g_j < 0\}$.

Repeat until either set Z or J is empty:

1. Repeat until either set Z or J is empty:

(a) Compute $\tilde{x}_J := \tilde{c}_J - C_{JJ}^{-1} \tilde{g}_J$.

(b) Update $\tilde{g}_Z := \tilde{g}_Z + A_{ZJ} \tilde{x}_J$.

(c) Redefine $J := \{j \in Z \mid g_j < 0\}$.

(d) Let $P := P \cup J$ and $Z := Z \setminus J$.

2. Compute $\tilde{x}_P := A_{PP}^{-1} (\tilde{b}_P - A_{PZ} \tilde{c}_Z)$.

3. Compute $\tilde{g}_Z := A_{ZP} \tilde{x}_P + A_{ZZ} \tilde{c}_Z - \tilde{b}_Z$.

4. Redefine $J := \{j \in Z \mid g_j < 0\}$.

5. Let $P := P \cup J$ and $Z := Z \setminus J$. \square

3.3. Lower and Upper Bounds Problem

In this section, we consider the lower and upper bounds quadratic program

$$\begin{aligned} \min_{\underline{x}} \quad & \frac{1}{2} \underline{x}^t A \underline{x} - \underline{x}^t \underline{b} \\ \text{subject to} \quad & \underline{c} \leq \underline{x} \leq \underline{d}. \end{aligned} \tag{1.1}$$

We introduce an index set U denoting the variables constrained at the upper bound \underline{d} . The index sets P and Z denote the unconstrained variables and the variables constrained at the lower bound \underline{c} , respectively.

Our lower and upper bounds algorithm starts with the iterate vector \underline{x} at the upper bound \underline{d} and the index set U equal to N . We examine the gradient vector $g(\underline{x}) \equiv A\underline{x} - \underline{b}$, and release those variables whose gradients are positive. Let P denote the just released variables and let $U := \bar{P}$. We solve the lower bound quadratic program

$$\begin{aligned} \min_{\underline{y}} \quad & \frac{1}{2} \underline{y}^t A_{PP} \underline{y} - \underline{y}^t (\underline{b}_P - A_{PU} \underline{d}_U) \\ \text{subject to} \quad & \underline{y} \geq \underline{c}_P. \end{aligned} \tag{3.1}$$

The iterate vector \underline{x} is redefined so that $\underline{x}_U = \underline{d}_U$ and $\underline{x}_P = \underline{y}$, the solution to (3.1). The indices of the components of \underline{x} at the lower bound are moved from set P to set Z . We start a new iteration by releasing those variables at the upper bound whose gradients have now become positive.

THEOREM 3.1. Suppose we have a vector \tilde{x} and two nonempty index sets U and P partitioning N such that $\tilde{x}_U = \underline{d}_U$ and $\tilde{x}_P = \tilde{y}$, where \tilde{y} solves (3.1). Let $J := \{j \in U | g_j(\tilde{x}) > 0\}$.

If the index set J is empty, then

1. \tilde{x} solves the lower and upper bounds quadratic program (1.1), else
2. let $P := \{j \in P | x_j > c_j\}$, $Z := \{j \in P | x_j = c_j\}$, $K := P \cup J$ and $L := U \setminus J$. Define a vector \tilde{y} so that $\tilde{y}_L = \underline{d}_L$, $\tilde{y}_Z = \underline{c}_Z$ and $\tilde{y}_K = p$, where p solves the lower bound quadratic program

$$\min_p \frac{1}{2} p^t A_{KK} p - p^t (b_K - A_{KL} \underline{d}_L)$$

subject to $p \geq \underline{c}_K$

Then

- (a) $\tilde{y}_J < \underline{d}_J$ and $\tilde{y}_P \leq \tilde{x}_P$.
- (b) $g_L(\tilde{y}) \geq g_L(\tilde{x})$ and $g_Z(\tilde{y}) \geq g_Z(\tilde{x}) \geq 0$.
- (c) $v(\tilde{y}) < v(\tilde{x})$.

Proof.

- (1) The optimality conditions (1.5) are satisfied because $g_U(\tilde{x}) \leq 0$ and \tilde{x}_P solves (3.1).
- (2a) The inequalities hold trivially for those components of \tilde{y}_K that are at the lower bound. The gradients of the other components equal zero and we can prove the inequalities using a technique similar to that of Part (2a) of Theorem 2.1.
- (2b) We have shown in Part (2a) that $\tilde{y}_K \leq \tilde{x}_K$. As $\tilde{y}_{\bar{K}} = \tilde{x}_{\bar{K}}$ by construction, we get $g_{\bar{K}}(\tilde{y}) \geq g_{\bar{K}}(\tilde{x})$ from Lemma 2.3. But $\bar{K} = L \cup Z$.

(2c) Since $g_Z(\underline{x}) \geq \underline{c}$ by Part (2b), the vector $\underline{g} = \underline{x}_K \cup Z$ solves the lower quadratic program

$$\min_{\underline{g}} \frac{1}{2} \underline{g}^t A_{QQ} \underline{g} - \underline{g}^t (\underline{b}_Q - A_{QL} \underline{d}_L)$$

subject to $\underline{g} \geq \underline{c}_Q$,

where $Q = K \cup Z$. But $\bar{L} = Q$, $\underline{x}_L = \underline{x}_L$ and $\underline{x} \neq \underline{x}$. Hence $v(\underline{x}) < v(\underline{x})$. \square

ALGORITHM 3.1. (Pang's method). Set $\underline{x} := \underline{d}$ and compute $\underline{g} := A\underline{d} - \underline{b}$. Let $J := \{j \in N \mid g_j > 0\}$, $P := J$, $U := \bar{J}$ and $Z := \emptyset$.

Repeat until set U or P or J is empty:

1. Reconstruct \underline{x} so that $\underline{x}_U = \underline{d}_U$, $\underline{x}_Z = \underline{c}_Z$ and $\underline{x}_P = \underline{x}$, where \underline{x} solves the lower bound quadratic program (3.1) (we may use Algorithm 2.3).
2. Set $K := \{k \in P \mid x_k = c_k\}$. Let $P := P \setminus K$ and $Z := Z \cup K$.
3. Compute $\underline{g}_U := A_{UU} \underline{d}_U + A_{UP} \underline{x}_P + A_{UZ} \underline{c}_Z - \underline{b}_U$.
4. Redefine $J := \{j \in U \mid g_j > 0\}$. Set $P := P \cup J$ and $U := U \setminus J$. \square

Part (2a) of Theorem 3.1 says that the iterate vector is non-increasing in value. Hence once a variable leaves its upper bound it never returns, and once a variable enters its lower bound it never exits. Since we release variables from their upper bounds when their gradients become positive, there is a flow of variables from the upper bound to the unconstrained region, and then to the lower bound. If

there is no outflow of variables from the upper bound, then the iterate vector is the desired solution (Part (1) of Theorem 3.1). Thus, Algorithm 3.1 always terminates. It is a descent method by Part (2c) of the same theorem. Algorithm 3.1 is essentially the algorithm proposed by Pang [8], except that he releases variables from their upper bounds when their gradients are nonnegative. The difference is so minor in real arithmetic that we expect both algorithms to produce identical results for almost all problems. Indeed, the two algorithms behaved identically in all our test examples.

We are interested in a "partial updating" technique that is similar to the one in the previous section. Such a technique may save many expensive solutions of lower bound quadratic programs. However, the presence of a lower bound restricts our choice of a "partial updating" matrix to a diagonal matrix. Recall that a positive definite matrix has a positive diagonal.

ASSUMPTION 3.1. We have chosen a diagonal matrix D such that $D \geq A$.

NOTATION. If $\underline{v} = (v_1, v_2, \dots, v_p)^t$, then $\underline{v}_+ = (w_1, w_2, \dots, w_p)^t$, where $w_j = \max(v_j, 0)$ for $j = 1, 2, \dots, p$.

LEMMA 3.2. Suppose there is a vector \underline{x} with $\underline{c} \leq \underline{x} \leq \underline{d}$ such that $\underline{x}_J = \underline{d}_J$ and $g_J(\underline{x}) > 0$ for some nonempty subset J of N . Construct a vector \underline{y} so that $\underline{y}_{\bar{J}} = \underline{x}_{\bar{J}}$ and $\underline{y}_J = \underline{c}_J + [\underline{d}_J - \underline{c}_J - D_{JJ}^{-1} g_J(\underline{x})]_+$. Then $g_J(\underline{y}) \geq 0$.

Proof. There exists a diagonal matrix $\hat{D} \geq D_{JJ}$ such that

$$x_J = c_J + [d_J - c_J - \hat{D}^{-1}g_J(x)] .$$

Hence $\hat{D} \geq A_{II}$ and this lemma can be proved in the same way as

Lemma 2.5. \square

Note that Lemma 3.2 may not hold if we replace D by a non-diagonal M-matrix E with $E \geq A$.

THEOREM 3.2. Suppose we have a vector x with $c < x \leq d$ and two nonempty index sets U and P partitioning N such that $x_U = d_U$, $x_P < d_P$ and $g_P(x) \geq 0$. Let $J := \{j \in U | g_j(x) > 0\}$ and assume that it is not empty. Construct a vector y such that $y_{\bar{J}} = x_{\bar{J}}$ and $y_J = c_J + [d_J - c_J - D_{JJ}^{-1}g_J(x)]_+$. Let $K := U \setminus J$. Then

- (1) $g_K(y) \geq g_K(x)$, $g_J(y) \geq 0$ and $g_P(y) \geq g_P(x) \geq 0$.
- (2) $y_P \geq q$, where q solves the lower bound quadratic program (3.1).
- (3) $v(y) < v(x)$.

Proof.

- (1) By Lemma 3.2, $g_J(y) \geq 0$. Since $y_J < d_J = x_J$ and $y_{\bar{J}} = x_{\bar{J}}$, we get $g_{\bar{J}}(y) \geq g_{\bar{J}}(x)$ by Lemma 2.3. But $\bar{J} = K \cup P$.
- (2) The inequality is trivial for those components of q that are at the lower bound. Since the gradients of the other components of q equal zero and $g_P(y) \geq 0$ from Part (1), we can complete the proof using a technique similar to the one in Lemma 2.6.

(3) Let $\hat{D} \geq D_{JJ}$ be the diagonal matrix such that

$$\begin{aligned} x_J &= c_J + [d_J - c_J - \hat{D}^{-1} g_J(x)] \\ &= d_J - \hat{D}^{-1} g_J(x). \end{aligned}$$

Let $z \equiv \hat{D}^{-1} g_J(x)$. Then

$$x_J = x_J - z$$

and

$$\begin{aligned} v(y) - v(x) &= \frac{1}{2} z^t A_{JJ} z - z^t A_{JJ} x_J + z^t b_J \\ &= \frac{1}{2} z^t A_{JJ} \hat{D}^{-1} g_J(x) - z^t g_J(x) \\ &= -z^t (I - \frac{1}{2} A_{JJ} \hat{D}^{-1}) g_J(x). \end{aligned}$$

It is easy to show that $I - (A_{JJ} \hat{D}^{-1})/2 \geq 0$, equality excluded.

Since $z > 0$ and $g_J(x) > 0$, we get $v(y) - v(x) < 0$. \square

ALGORITHM 3.2. (Modified Pang's method). Set $x := d$ and compute $g := Ad - b$.

Let $J := \{j \in N | g_j > 0\}$, $P := J$, $U := \bar{J}$ and $Z := \emptyset$.

Repeat until set U or P or J is empty:

1. Repeat until set U or J is empty:

(a) Update $x_J := c_J + [d_J - c_J - D_{JJ}^{-1} g_J]_+$.

(b) Let $K := \{k \in J | x_k = c_k\}$, $Z := Z \cup K$ and $P := P \cup (J \setminus K)$.

(c) Update $g_U := g_U - A_{UU}(d_J - x_J)$.

(d) Redefine $J := \{j \in U | g_j > 0\}$. Let $U := U \setminus J$.

2. Reconstruct \tilde{x} so that $\tilde{x}_U = \tilde{d}_U$, $\tilde{x}_Z = \tilde{c}_Z$ and $\tilde{x}_P = y$, where y solves the lower bound quadratic program (3.1) (we may use Algorithm 2.3).
3. Set $K := \{k \in P \mid x_k = c_k\}$. Let $P := P \setminus K$ and $Z := Z \cup K$.
4. Compute $g_U := A_{UU}d_U + A_{UP}x_P + A_{UZ}c_Z - b_U$.
5. Redefine $J := \{j \in U \mid g_j > 0\}$. Let $P := P \cup J$ and $U := U \setminus J$. \square

3.4. Problem with Non-finite Bounds

In this section, we consider the lower and upper bounds quadratic program

$$\begin{aligned} \min_{\tilde{x}} \quad & \frac{1}{2} \tilde{x}^t A \tilde{x} - \tilde{x}^t b \\ \text{subject to} \quad & \tilde{c} \leq \tilde{x} \leq \tilde{d}, \end{aligned} \tag{1.1}$$

where some components of the bounds are infinite. We show that we can compute finite bounds on all the variables.

LEMMA 4.1. Consider the quadratic program (1.1) with $\tilde{d} = \infty$. Assume $\tilde{c}_J = -\infty$ and $\tilde{c}_{\bar{J}} > -\infty$ for some nonempty subset J of N . Let \tilde{x} be the solution vector. Then

$$\tilde{x}_J \geq A_{JJ}^{-1}(b_J - A_{J\bar{J}}\tilde{c}_{\bar{J}}).$$

Proof. Define a vector p with $p_J = A_{JJ}^{-1}(b_J - A_{J\bar{J}}\tilde{c}_{\bar{J}})$ and $p_{\bar{J}} = \tilde{c}_{\bar{J}}$. By construction, $g_J(p) = 0$. Define another vector q with $q_J = p_J$ and $q_{\bar{J}} = \tilde{x}_{\bar{J}}$. Since $q_{\bar{J}} \geq p_{\bar{J}}$, we get $g_J(q) \leq g_J(p)$ from Lemma 2.3.

Also $g_J(\tilde{x}) = 0$ by the optimality conditions (1.5). Thus,

$$\begin{aligned} A_{JJ}(\tilde{x}_J - p_J) &= A_{JJ}\tilde{x}_J + A_{J\bar{J}}\tilde{x}_{\bar{J}} - b_J - (A_{JJ}p_J + A_{J\bar{J}}\tilde{x}_{\bar{J}} - b_J) \\ &= g_J(\tilde{x}) - g_J(p) \\ &\geq 0. \end{aligned}$$

Hence

$$\tilde{x}_J - p_J \geq 0. \quad \square$$

LEMMA 4.2. Consider the quadratic program (1.1) with $c = -\infty$. Assume $d_J = \infty$ and $d_{\bar{J}} < \infty$ for some nonempty subset J of N . Let \tilde{x} be the solution vector. Then

$$\tilde{x}_J \leq A_{JJ}^{-1}(b_J - A_{J\bar{J}}\tilde{x}_{\bar{J}}).$$

Proof. Similar to that of Lemma 4.1. \square

We now have the tools to handle the problem when both bounds have infinite components.

LEMMA 4.3. For problem (1.1), assume $d_J = \infty$ and $d_{\bar{J}} < \infty$ for some nonempty subset J of N . Let \tilde{x} be the solution to (1.1) and y be the solution to the lower bound quadratic program

$$\begin{aligned} \min_y \quad & \frac{1}{2} y^t A_{JJ} y - y^t (b_J - A_{J\bar{J}} d_{\bar{J}}) \\ \text{subject to} \quad & y \geq c_J. \end{aligned} \tag{4.1}$$

Then

$$\tilde{x}_J \leq y.$$

Proof. The inequality is trivial for those components of \tilde{x}_J that are at the lower bound. Let $K := \{k \in J \mid \tilde{x}_k > c_k\}$. Then $g_K(\tilde{x}) = 0$.

Now define a vector p with $p_J = \tilde{x}$ and $p_{\bar{J}} = \underline{d}_{\bar{J}}$. Hence $g_J(p) \geq 0$ by the optimality conditions (1.5). Define another vector q with $q_J = \tilde{x}$ and $q_{\bar{J}} = \tilde{x}_{\bar{J}}$. Since $q_{\bar{J}} \leq p_{\bar{J}}$, we get $g_J(q) \geq g_J(p)$ by Lemma 2.3. It follows that $g_K(q) \geq 0$.

Let $L := J \setminus K$. Then $\tilde{x}_L = c_L \leq p_L$, or equivalently, $\tilde{x}_L - p_L \leq 0$. Thus,

$$\begin{aligned} A_{KK}(\tilde{x}_K - p_K) &= (A_{KK}\tilde{x}_K + A_{KL}\tilde{x}_L + A_{K\bar{J}}\tilde{x}_{\bar{J}} - b_K) \\ &\quad - (A_{KK}p_K + A_{KL}p_L + A_{K\bar{J}}\tilde{x}_{\bar{J}} - b_K) - A_{KL}(\tilde{x}_L - p_L) \\ &= g_K(\tilde{x}) - g_K(q) - A_{KL}(\tilde{x}_L - p_L) \\ &\leq 0 \quad \text{because } A_{KL} \leq 0. \end{aligned}$$

Hence

$$\tilde{x}_K - p_K \leq 0. \quad \square$$

LEMMA 4.4. For problem (1.1), assume $c_J = -\infty$ and $c_{\bar{J}} > -\infty$ for some nonempty subset J of N . Let \tilde{x} be the solution to (1.1) and \underline{x} be the solution to the upper bound quadratic program

$$\begin{aligned} \min_{\underline{x}} \quad & \frac{1}{2} \underline{x}^t A_{JJ} \underline{x} - \underline{x}^t (b_J - A_{J\bar{J}} c_{\bar{J}}) \\ \text{subject to} \quad & \underline{x} \leq \underline{d}_J. \end{aligned} \tag{4.2}$$

Then

$$\tilde{x}_J \geq \underline{x}.$$

Proof. Similar to that of Lemma 4.3. \square

We can solve the lower bound quadratic program (4.1) by first using Lemma 4.1 to compute finite lower bounds for those variables with

constraints and then applying Algorithm 2.3. The upper bound problem (4.2) can be solved in a similar fashion using Lemma 4.2 and Algorithm 2.3. Note that an upper bound problem is reduced to one with a lower bound if all its variables are negated.

3.5. Numerical Examples

We have chosen four representative problems to study the effectiveness of our "partial updating" technique. We use the conjugate gradient method as our matrix equations solver. Our programs were written in FORTRAN and run on an IBM 370/168 computer at the Stanford Linear Accelerator Center. The codes were compiled with optimization level 2 of the H EXTENDED compiler.

EXAMPLE 1. Let us consider the linear complementarity problem

$$\begin{aligned} \min_{\tilde{x}} \quad & \frac{1}{2} \tilde{x}^t A \tilde{x} - \tilde{x}^t b \\ \text{subject to} \quad & \tilde{x} \geq 0. \end{aligned} \tag{1.2}$$

The matrix A is chosen as

$$A = \begin{pmatrix} 2 & -1 & & & \bigcirc \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ \bigcirc & -1 & 2 & -1 & \\ & \bigcirc & -1 & 2 & \end{pmatrix}_{n \times n}$$

and the vector \underline{b} is generated by

$$b_i = 8 - 20r_i \quad \text{for } i = 1, 2, \dots, n,$$

where r_i is a random number from a uniform distribution in the open interval $(0,1)$. The "partial updating" matrix C is chosen as

$$C = \begin{pmatrix} 2 & & & \bigcirc \\ & 2 & & \\ & & \ddots & \\ \bigcirc & & & 2 \end{pmatrix}_{n \times n}.$$

We define an iteration to be a sequence of "partial updates" followed by a solution of a matrix equation with coefficient matrix A_{PP} . The scalar $|P|$ gives the number of elements in the index set P at the end of an iteration. The other scalar u gives the number of "partial updates" in an iteration. Time is machine execution time in seconds.

n = 1000

Iteration	Algorithm 2.1	Algorithm 2.2
	P	u, P
1	382	2,456
2	499	2,524
3	535	1,539
4	539	1,542
5	542	
Time	2.67	2.40

n = 1500

Iteration	Algorithm 2.1	Algorithm 2.2
	P	u, P
1	585	2,692
2	763	2,806
3	824	3,838
4	837	2,848
5	845	2,851
6	848	1,853
7	850	
8	852	
9	853	
Time	8.43	6.49

n = 2000

Iteration	Algorithm 2.1	Algorithm 2.2
	P	u, P
1	782	2,923
2	1020	2,1074
3	1098	3,1118
4	1119	2,1131
5	1128	2,1136
6	1133	1,1138
7	1135	
8	1137	
9	1138	
Time	11.79	9.92

AD-A065 285

STANFORD UNIV CALIF DEPT OF COMPUTER SCIENCE
SPARSE AND PARALLEL MATRIX COMPUTATIONS.(U)

F/G 12/1

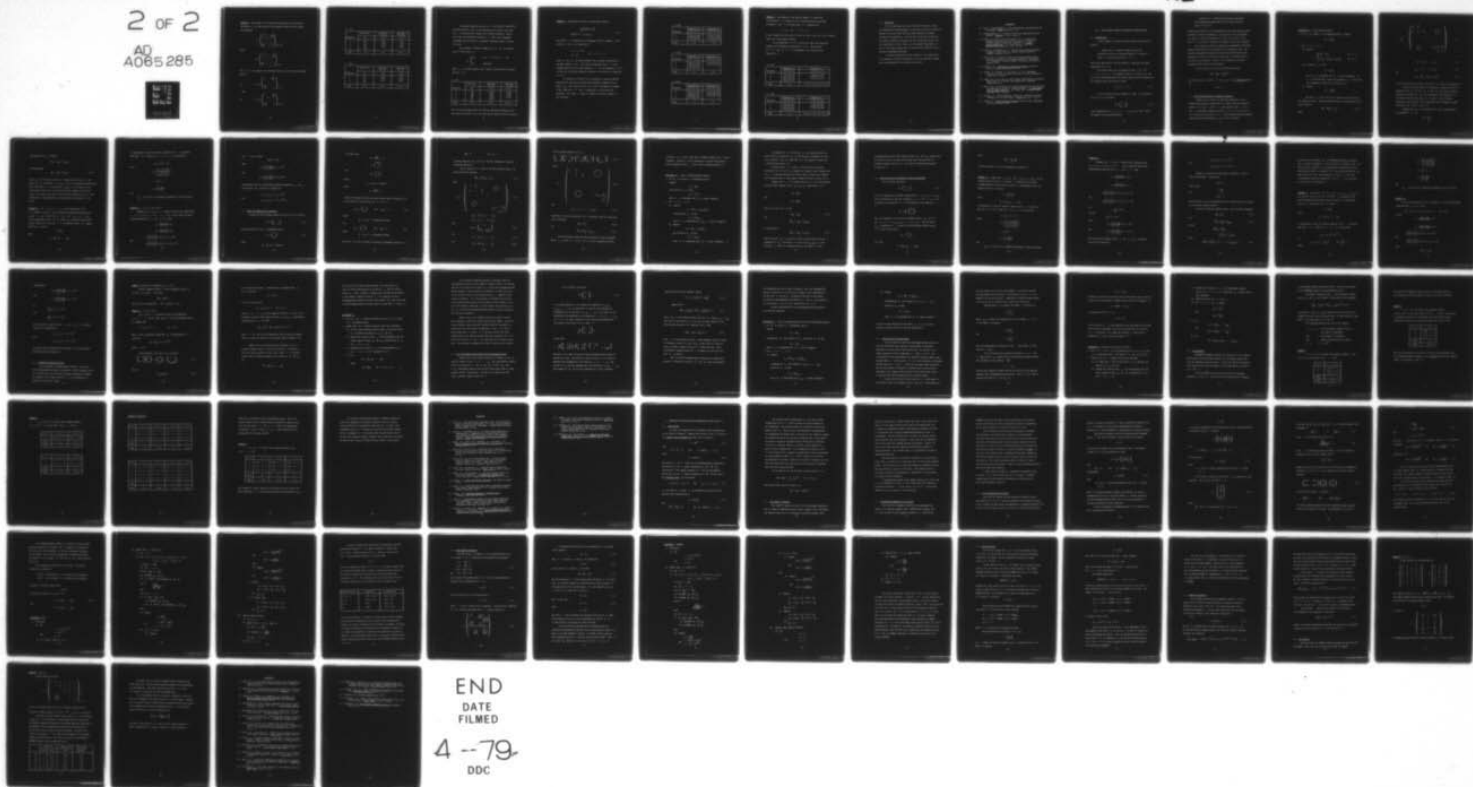
DEC 78 F T LUK
STAN-CS-78-685

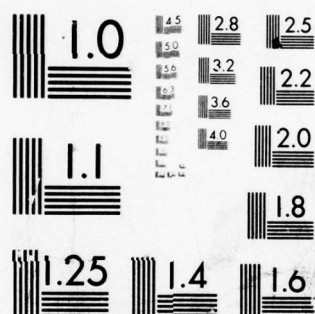
DAHC04-75-6-0185
NL

UNCLASSIFIED

2 OF 2

AD
A065285





MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

EXAMPLE 2. This example is the same as the previous one, except that the matrix A has been chosen as the Laplacian 5-point finite difference operator:

$$A = \begin{pmatrix} B & -I & & \bigcirc \\ -I & B & -I & \\ \cdot & \cdot & \cdot & \cdot \\ \bigcirc & -I & B & -I \\ & -I & B & \end{pmatrix}_{m^2 \times m^2},$$

where

$$B = \begin{pmatrix} 4 & -1 & & \bigcirc \\ -1 & 4 & -1 & \\ \cdot & \cdot & \cdot & \cdot \\ \bigcirc & -1 & 4 & -1 \\ & -1 & 4 & \end{pmatrix}_{m \times m}.$$

Let $n = m^2$. We consider two different choices of the "partial updating" matrix:

$$(1) \quad C = D \equiv \begin{pmatrix} 4 & & \bigcirc \\ & 4 & \\ \bigcirc & \cdot & \cdot \\ & & 4 \end{pmatrix}_{n \times n},$$

and

$$(2) \quad C = T \equiv \begin{pmatrix} B & & \bigcirc \\ & B & \\ \bigcirc & \cdot & \cdot \\ & & B \end{pmatrix}_{n \times n}.$$

n = 900

	Algorithm 2.1	Algorithm 2.2 (C = D)	Algorithm 2.2 (C = T)
Iteration	P	u, P	u, P
1	349	2,416	2,416
2	450	2,462	2,462
3	463	1,465	1,465
4	465	1,466	1,466
5	466		
Time	7.54	6.73	6.70

n = 1600

	Algorithm 2.1	Algorithm 2.2 (C = D)	Algorithm 2.2 (C = T)
Iteration	P	u, P	u, P
1	620	2,730	2,730
2	798	2,828	1,828
3	837	1,842	1,842
4	842		
Time	17.08	15.29	15.41

We observe that both choices of C have produced essentially identical results. This is not surprising if we look at the index set J of the variables eligible for "partial updating." Rarely do we find two consecutive indices in J . Thus, the tridiagonal "updating" matrix reduces to a diagonal "updating" matrix in almost all cases.

We construct a different example with $n = 900$ and another choice of the vector b :

$$b_i = \begin{cases} -6r_i & \text{for } i = 301, 302, \dots, 600 \\ 8r_i & \text{otherwise,} \end{cases}$$

where r_i is a random number from a uniform distribution in the open interval $(0,1)$.

$n = 900$

	Algorithm 2.1	Algorithm 2.2 ($C = D$)	Algorithm 2.2 ($C = T$)
Iteration	$ P $	$u, P $	$u, P $
1	300	2,307	2,307
2	366	2,380	3,386
3	396	2,401	2,405
4	406	1,406	1,407
5	407	1,407	
Time	24.97	25.08	23.80

This is also an example where the "partial updating" technique is not particularly efficient due to the very special structure of the vector b .

EXAMPLE 3. We address the lower and upper bounds problem

$$\begin{aligned} \min_{\tilde{x}} \quad & \frac{1}{2} \tilde{x}^t A \tilde{x} - \tilde{x}^t b \\ \text{subject to} \quad & 0 \leq \tilde{x} \leq d. \end{aligned} \tag{5.1}$$

The matrix A is chosen as the tridiagonal matrix in Example 1. The vectors b and d are generated by

$$\begin{aligned} b_i &= 11 - 20r_i, \\ d_i &= 7s_i, \end{aligned} \quad \text{for } i = 1, 2, \dots, n,$$

where r_i and s_i are random numbers from a uniform distribution in the open interval $(0,1)$. The "partial updating" matrix C is the diagonal matrix with the same diagonal as A . We use Algorithm 2.1 (2.2) to solve the lower bound quadratic program in an iteration of Algorithm 3.1 (3.2).

An iteration is defined to be a sequence of "partial updates" followed by a solution of a lower bound quadratic program with the matrix A_{pp} . The 3-tuple $(|Z|, |P|, |U|)$ gives the number of elements in the index sets Z , P and U , respectively, at the end of an iteration. The scalar u gives the number of "partial updates" in the iteration.

n = 1000

Iteration	Algorithm 3.1	Algorithm 3.2
	(Z , P , U)	u, (Z , P , U)
1	(234, 216, 550)	3, (307, 275, 418)
2	(310, 283, 407)	1, (311, 292, 397)
3	(313, 291, 396)	1, (313, 291, 396)
Time	1.80	1.63

n = 1500

Iteration	Algorithm 3.1	Algorithm 3.2
	(Z , P , U)	u, (Z , P , U)
1	(366, 324, 810)	3, (489, 396, 615)
2	(492, 408, 600)	1, (495, 416, 589)
3	(495, 415, 590)	
4	(495, 416, 589)	
Time	4.13	2.09

n = 2000

Iteration	Algorithm 3.1	Algorithm 3.2
	(Z , P , U)	u, (Z , P , U)
1	(452, 452, 1096)	3, (590, 560, 850)
2	(602, 579, 819)	2, (605, 602, 793)
3	(605, 602, 793)	1, (605, 603, 792)
4	(605, 603, 792)	
Time	7.51	5.01

EXAMPLE 4. This example is the same as Example 3, except that

- (1) the matrix A is chosen as the discretized Laplacian operator of Example 2, and (2) the upper bound \underline{d} is generated by

$$d_i = 3s_i \quad \text{for } i = 1, 2, \dots, n.$$

We have changed the upper bound so that the three index sets are of roughly equal size in the final solution.

We use Algorithm 2.1 (2.2) to solve the lower bound quadratic program in an iteration of Algorithm 3.1 (3.2). Both choices:

- (1) $C = D$ and (2) $C = T$ (see Example 2) are considered for Algorithm 2.2.

$n = 900$

Iteration	Algorithm 3.1	Algorithm 3.2
	(Z , P , U)	$u, (Z , P , U)$
1	(274, 168, 458)	3, (307, 223, 370)
2	(325, 277, 298)	2, (325, 294, 281)
3	(325, 294, 281)	1, (325, 296, 279)
4	(325, 296, 279)	
Time	6.79	5.97 ($C = D$), 5.94 ($C = T$)

$n = 1600$

Iteration	Algorithm 3.1	Algorithm 3.2
	(Z , P , U)	$u, (Z , P , U)$
1	(478, 277, 845)	3, (552, 430, 618)
2	(582, 487, 531)	2, (585, 524, 497)
3	(585, 520, 495)	1, (585, 526, 489)
4	(585, 526, 489)	
Time	17.72	15.35 ($C = D$), 15.28 ($C = T$)

3.6. Conclusion

It is evident from the work of Cottle and Goheen [4] that the algorithms of Chandrasekaran [2] and Pang [8] are very effective schemes for solving lower and upper bounds quadratic programs associated with M-matrices. We have seen in the last section how our "partial updating" technique can cut the execution time of the two algorithms by 10-50%. Thus, our new schemes (Algorithms 2.2 and 3.2) are highly competitive for solving this important class of quadratic programming problems.

Our other contribution of this chapter is the introduction of a technique to handle the special case of the quadratic program when some components of the bounds are infinite.

REFERENCES

- [1] Ceá, J., and Glowinski, R., "Sur des methodes d'optimisation par relaxation," RAIRO R-3 (1973), 5-32.
- [2] Chandrasekaran, R., "A special case of the complementary pivot problem," Opsearch 7 (1970), 263-268.
- [3] Concus, P., Golub, G.H., and O'Leary, D.P., "A generalized conjugate gradient method for the numerical solution of elliptic partial differential equations," in Sparse Matrix Computations (J. R. Bunch and D.J. Rose, Ed.) Academic Press, New York (1976), 309-332.
- [4] Cottle, R. and Goheen, M., "A special class of large quadratic programs," Report SOL-76-7, Systems Optimization Lab., Stanford University (1976).
- [5] Cryer, C.W., "A survey of trial free-boundary methods for the numerical solution of free-boundary problems," MRC Report 1693, Mathematics Research Center, University of Wisconsin, Madison (1976).
- [6] Fröberg, C.E., Introduction to Numerical Analysis, Addison-Wesley, Reading, Massachusetts (1965).
- [7] Levati, G., Scarpini, F., and Volpi, G., "Sul trattamento numerico di alcuni problemi variazionali di tipo unilaterale," L.A.N. Pub. 82 (1974).
- [8] Pang, J.S., "On a class of least-element complementarity problems," Report SOL-76-10, Systems Optimization Lab., Stanford University (1976).
- [9] Reid, J.K., "On the method of conjugate gradients for the solution of large sparse systems of linear equations," in Large Sparse Sets of Linear Equations, (J.K. Reid, Ed.), Academic Press, New York (1971), 231-254.
- [10] Scarpini, F., "Some algorithms solving the unilateral Dirichlet problems with two constraints," Calcolo 12 (1975), 113-149.
- [11] Varga, R.S., Matrix Iterative Analysis, Prentice Hall, Englewood Cliffs, New Jersey (1962).

IV. A BLOCK LANCZOS METHOD FOR COMPUTING SINGULAR VALUES

4.1. Introduction

In this chapter, we construct a Block Lanczos method for the problem:

Compute the k greatest singular values and associated vectors of a large and sparse $m \times n$ matrix A , where k is much smaller than m or n ,

which finds applications in factor analysis, regression and image enhancement (cf. [4]).

We assume without loss of generality that $m \geq n$. For $i = 1, 2, \dots, n$, let σ_i be a singular value of A , and let \underline{u}_i and \underline{v}_i be the corresponding left and right singular vectors, respectively. The singular values are ordered so that

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n. \quad (1.1)$$

Let us exploit an idea of Lanczos [8, Chap. 3] and consider the $(m+n) \times (m+n)$ matrix

$$\tilde{A} = \begin{pmatrix} 0 & A \\ A^t & 0 \end{pmatrix}, \quad (1.2)$$

whose eigenvalues are $\pm \sigma_1, \pm \sigma_2, \dots, \pm \sigma_n$, plus $(m-n)$ zeros.

We address the equivalent problem:

Compute the k algebraically greatest eigenvalues and corresponding eigenvectors of the large and sparse matrix \tilde{A} of (1.2).

An efficient scheme for this eigenproblem is the Block Lanczos method developed by several researchers, in particular, Cullum and Donath ([1] and [2]), Golub and Underwood ([7] and [14]), Lewis [9] and Ruhe [11]. We choose to consider the variant of Golub and Underwood.

We are going to present a theoretical development of the Block Lanczos method and give two theorems on its convergence rate. The practical implementation aspects are then discussed and particular attention is paid to the choice of block size. We believe that all our results are original. In fact, we are unaware of any other procedures which solve the same problem.

In this chapter, we use the Euclidean vector norm

$$\|\tilde{x}\| = \|\tilde{x}\|_2 = (\tilde{x}^t \tilde{x})^{1/2},$$

and refer to an $n \times b$ matrix X with $n \geq b$ as an orthonormal matrix if

$$X^t X = I.$$

4.2. Block Lanczos Method for Symmetric Matrices

Suppose that we desire accurate approximations to the k algebraically greatest eigenvalues of a large, sparse and symmetric matrix B of order ℓ , where k is much smaller than ℓ . Let b and s be two given integers such that $b \geq 1$, $s \geq 2$ and $bs \leq \ell$. It is usually the case that $bs \ll \ell$. We can define the Block Lanczos method of Golub and Underwood ([7] and [14]) as follows.

ALGORITHM 2.1. (Block Lanczos method)

1. Let X_1 be a given $\ell \times b$ orthonormal matrix. Compute

$$M_1 := X_1^t B X_1 .$$

2. For $i = 2, 3, \dots, s$, do

- (a) Compute

$$Z_i := \begin{cases} B X_1 - X_1 M_1 & \text{for } i = 2, \\ B X_{i-1} - X_{i-1} M_{i-1} - X_{i-2} R_{i-1}^t & \text{for } i \geq 3. \end{cases}$$

- (b) Factorize Z_i so that

$$Z_i := X_i R_i ,$$

where X_i is orthonormal and R_i is upper triangular. If Z_i is rank deficient, choose the columns of X_i so that they are orthogonal to those of all previous X_j 's.

- (c) Compute

$$M_i := X_i^t B X_i .$$

□

The value b is thus the order of the block, and the value s the number of blocks. The Block Lanczos method is characterized by the matrix equation

$$B \bar{X}_s = \bar{X}_s \bar{M}_s + \bar{Z}_{s+1} , \quad (2.1)$$

where

$$\bar{M}_s = \begin{pmatrix} M_1 & R_2^t & & & \\ R_2 & M_2 & R_3^t & & \\ & \cdot & \cdot & \cdot & \\ & & \cdot & \cdot & \cdot \\ & & & \cdot & \cdot \\ & & & & \cdot \\ & & & & & R_{s-1}^t & M_{s-1}^t & R_s^t \\ & & & & & & R_s & M_s \end{pmatrix}, \quad (2.2)$$

$$\bar{X}_s = (X_1, X_2, \dots, X_s), \quad (2.3)$$

$$\bar{Z}_{s+1} = (0, \dots, 0, Z_{s+1}) \quad (2.4)$$

and

$$Z_{s+1} = BX_s - X_s M_s - X_{s-1} R_s^t. \quad (2.5)$$

We can prove that the columns of the matrix \bar{X}_s form an orthonormal set.

The Block Lanczos method thus generates a symmetric block tri-diagonal matrix \bar{M}_s of order bs . As the R_i 's are upper triangular matrices, the matrix \bar{M}_s is also a band matrix with bandwidth $2b + 1$. We consider only the case where bs is small, so that standard techniques can be applied to the computation of the eigenvalue decomposition of \bar{M}_s (see [12] and [15]).

Suppose that y_i is an eigenvector of \bar{M}_s corresponding to the eigenvalue μ_i . Let

$$x_i = \bar{X}_s y_i.$$

From equation (2.1), we obtain

$$B\tilde{x}_i = \mu_i \tilde{x}_i + \tilde{z}_{s+1} y_i ,$$

or equivalently,

$$B\tilde{x}_i = \mu_i \tilde{x}_i + Z_{s+1} f_i , \quad (2.6)$$

where f_i is the vector of order b consisting of the last b components of y_i . Therefore, μ_i and \tilde{x}_i would be an eigenvalue-eigenvector pair of the matrix B if we had that $Z_{s+1} f_i = 0$. The power of the Block Lanczos method lies in the fact that the elements of the vector $Z_{s+1} f_i$ are usually very small for the extreme eigenvalues of \tilde{M}_s . This observation is substantiated by a theorem given by Underwood [14, pp. 37-38].

THEOREM 2.1. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_\ell$ be the eigenvalues of an $\ell \times \ell$ symmetric matrix B with corresponding normalized eigenvectors g_1, g_2, \dots, g_ℓ . Assume that $\lambda_b > \lambda_{b+1}$. Let $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{bs}$ be the eigenvalues of the $bs \times bs$ matrix \tilde{M}_s generated by the Block Lanczos method with an initial $\ell \times b$ orthonormal matrix X_1 . Suppose that the $b \times b$ matrix

$$W = G_1^t X_1 ,$$

where

$$G_1 = (g_1, g_2, \dots, g_b) ,$$

is nonsingular so that its smallest singular value τ is positive.

Note that $\tau \leq 1$. Then, for $i = 1, 2, \dots, b$, we have that

$$\lambda_i \geq \mu_i \geq \lambda_i - \epsilon_i^2,$$

where

$$\epsilon_i^2 = \frac{(\lambda_i - \lambda_\ell) \tan^2 \theta}{T_{s-1}^2 \left(\frac{1 + r_i}{1 - r_i} \right)},$$

$$\theta = \cos^{-1} \tau,$$

$$r_i = \frac{\lambda_i - \lambda_{b+1}}{\lambda_i - \lambda_\ell}$$

and

T_{s-1} is the $(s-1)$ -st Chebyshev polynomial of the first kind.

EXAMPLE 2.1. (cf. [14, pp. 43-44])

Suppose that B is an $\ell \times \ell$ symmetric matrix with eigenvalues $\lambda_1 = 1.0$, $\lambda_2 = 0.9$, $\lambda_3 = 0.5$, \dots , $\lambda_\ell = 0.0$. Let us apply the Block Lanczos method with $b = 2$ and $s = 10$. Then

$$r_1 = \frac{1.0 - 0.5}{1.0 - 0.0} = \frac{1}{2},$$

$$r_2 = \frac{0.9 - 0.5}{0.9 - 0.0} = \frac{4}{9},$$

$$T_9 \left(\frac{1 + r_1}{1 - r_1} \right) = T_9(3) \doteq 3.9 \times 10^6$$

and

$$T_9 \left(\frac{1 + r_2}{1 - r_2} \right) = T_9(2.6) \doteq 9.8 \times 10^5.$$

Let $\tau = 0.04$, so that

$$\tan^2 \theta = 624 .$$

Now,

$$\epsilon_1^2 \doteq \frac{1.0 \times 624}{1.5 \times 10^{13}} \doteq 4.1 \times 10^{-11}$$

and

$$\epsilon_2^2 \doteq \frac{0.9 \times 624}{9.6 \times 10^{11}} \doteq 5.8 \times 10^{-10} .$$

Consequently, the two algebraically greatest eigenvalues μ_1 and μ_2 of the matrix \bar{M}_s satisfy the inequalities

$$\lambda_1 \geq \mu_1 \geq \lambda_1 - 4.1 \times 10^{-11}$$

and

$$\lambda_2 \geq \mu_2 \geq \lambda_2 - 5.8 \times 10^{-10} .$$

4.3. Block Bidiagonalization Algorithm

We apply the Block Lanczos method to the $(m+n) \times (m+n)$ matrix

$$\tilde{A} = \begin{pmatrix} 0 & A \\ A^t & 0 \end{pmatrix} \quad (1.2)$$

with the initial $(m+n) \times b$ orthonormal matrix

$$X_1 = \begin{pmatrix} 0 \\ Q_1 \end{pmatrix}, \quad (3.1)$$

where

Q_1 is an $n \times b$ matrix.

It follows that

$$M_1 = X_1^t \tilde{A} X_1 = 0 ,$$

$$Z_1 = \begin{pmatrix} A Q_1 \\ 0 \end{pmatrix}$$

and

$$X_2 = \begin{pmatrix} P_1 \\ 0 \end{pmatrix},$$

where

P_1 is an $m \times b$ matrix.

Thus,

$$M_2 = X_2^t \tilde{A} X_2 = 0 .$$

Using the relations defining the Block Lanczos method (Algorithm 2.1), we can prove by induction that, for $j = 1, 2, \dots$,

$$X_{2j-1} = \begin{pmatrix} 0 \\ Q_j \end{pmatrix} \quad \text{and} \quad M_{2j-1} = 0 , \quad (3.2)$$

where

Q_j is an $n \times b$ orthonormal matrix,

and

$$X_{2j} = \begin{pmatrix} P_j \\ 0 \end{pmatrix} \quad \text{and} \quad M_{2j} = 0 , \quad (3.3)$$

where

P_j is an $m \times b$ orthonormal matrix.

Since the X_j 's form a sequence of mutually orthonormal matrices, i.e.

$$X_i^t X_j = 0 \quad \text{for } i \neq j ,$$

we deduce that the P_j 's and Q_j 's form two sequences of mutually orthonormal matrices.

Let us carry out $2s$ steps of the Block Lanczos scheme. We obtain the matrix equation

$$\tilde{A}\tilde{X}_{2s} = \tilde{X}_{2s} T_{2s} + \tilde{Z}_{2s} , \quad (3.4)$$

where

$$T_{2s} = \begin{pmatrix} 0 & R_2^t & & & \bigcirc \\ R_2 & 0 & R_3^t & & \\ & \cdot & \cdot & \cdot & \\ & & \cdot & \cdot & \cdot \\ \bigcirc & & & \cdot & \cdot \\ & & & & R_{2s-1}^t & 0 & R_{2s}^t \\ & & & & R_{2s} & 0 \end{pmatrix} , \quad (3.5)$$

$$\tilde{X}_{2s} = (X_1, X_2, \dots, X_{2s}),$$

$$\tilde{Z}_{2s} = (0, \dots, 0, Z_{2s+1})$$

and

$$Z_{2s+1} = \begin{pmatrix} 0 \\ A^t P_s - Q_s^t R_{2s} \end{pmatrix} . \quad (3.6)$$

Let

$$\bar{P}_s = (P_1, P_2, \dots, P_s) \quad (3.7)$$

and

$$\bar{Q}_s = (Q_1, Q_2, \dots, Q_s) . \quad (3.8)$$

We can rewrite equation (3.4) as

$$\begin{pmatrix} 0 & A \\ A^t & 0 \end{pmatrix} \begin{pmatrix} \bar{P}_s & 0 \\ 0 & \bar{Q}_s \end{pmatrix} = \begin{pmatrix} \bar{P}_s & 0 \\ 0 & \bar{Q}_s \end{pmatrix} \begin{pmatrix} 0 & J_s \\ J_s^t & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ \hat{Z}_{2s+1} & 0 \end{pmatrix}, \quad (3.9)$$

where

$$J_s = \begin{pmatrix} R_2 & R_3^t & & & & \\ & R_4 & R_5^t & & & \\ & & \cdot & \cdot & & \\ & & & \cdot & \cdot & \\ & & & & \cdot & \\ & & & & & \cdot & \\ & & & & & & R_{2s-2}^t & R_{2s-1}^t \\ & & & & & & & R_{2s}^t \\ & & & & & & & & R_{2s} \end{pmatrix} \quad (3.10)$$

and

$$\hat{Z}_{2s+1} = (0, \dots, 0, Z_{2s+1}). \quad (3.11)$$

Furthermore, the matrix equation (3.9) is reducible into two lower-order matrix equations:

$$A\bar{Q}_s = \bar{P}_s J_s \quad (3.12)$$

and

$$A^t \bar{P}_s = \bar{Q}_s J_s^t + \hat{Z}_{2s+1}. \quad (3.13)$$

The Block Lanczos method therefore generates a block bidiagonal matrix J_s of order bs . As the R_i 's are upper triangular matrices,

the matrix J_s is also a band upper triangular matrix with b super-diagonals. Therefore, we have constructed a scheme that produces a block bidiagonal matrix J_s from a given rectangular matrix A .

ALGORITHM 3.1. (Block Bidiagonalization method).

1. Let Q_1 be a given $n \times b$ orthonormal matrix.

Compute

$$W_1 := AQ_1,$$

and factorize W_1 so that

$$W_1 := P_1 R_2,$$

where P_1 is orthonormal and R_2 is upper triangular.

2. For $i = 2, 3, \dots, s$, do

(a) Compute

$$Z_i := A^t P_{i-1} - Q_{i-1} R_{2i-2}^t,$$

and factorize Z_i so that

$$Z_i := Q_i R_{2i-1},$$

where Q_i is orthonormal and R_{2i-1} is upper triangular.

(b) Compute

$$W_i := AQ_i - P_{i-1} R_{2i-1}^t,$$

and factorize W_i so that

$$W_i := P_i R_{2i},$$

where P_i is orthonormal and R_{2i} is upper triangular. \square

In Algorithm 3.1, if the matrix Z_i were rank deficient, we would choose the columns of Q_i so that they are orthogonal to those of all previous Q_j 's (cf. Algorithm 2.1). The remedy is similar for a rank-deficient matrix W_i .

We assume that bs is small, so that we can use standard techniques (see [3] and [6]) to compute the singular value decomposition of J_s . A thorough discussion of various ways to compute the singular value decomposition of a band upper triangular matrix is given in [5].

Suppose that μ_i is a singular value of J_s with corresponding left and right singular vectors \tilde{w}_i and \tilde{z}_i , respectively. Let

$$p_i = \bar{P}_s \tilde{w}_i$$

and

$$g_i = \bar{Q}_s \tilde{z}_i.$$

From (3.12) and (3.13), we get

$$A g_i = \mu_i p_i \tag{3.14}$$

and

$$A^t p_i = \mu_i g_i + \hat{Z}_{2s+1} \tilde{w}_i,$$

or equivalently,

$$A^t p_i = \mu_i g_i + Z_{2s+1} h_i, \tag{3.15}$$

where the vector h_i is a vector of order b consisting of the last b components of \tilde{w}_i . Accordingly, if we had that $Z_{2i+1} h_i = 0$, then the value μ_i would be a singular value of the matrix A with

corresponding left and right singular vectors p_i and q_i , respectively. In the next section, we give error bounds which indicate that the greatest singular values of J_s are usually accurate approximations to those of A .

4.4. Error Bounds For the Singular Value Approximations

Let us consider the matrix

$$\tilde{A} = \begin{pmatrix} 0 & A \\ A^t & 0 \end{pmatrix}. \quad (1.2)$$

Its $(b+1)$ algebraically greatest eigenvalues are $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{b+1}$, and its algebraically smallest eigenvalue is $-\sigma_1$. For $i = 1, 2, \dots, b$, the normalized eigenvector of \tilde{A} corresponding to the eigenvalue σ_i is

$$g_i = \frac{1}{\sqrt{2}} \begin{pmatrix} u_i \\ v_i \end{pmatrix}.$$

Now, the eigenvalues of the block tridiagonal matrix T_{2s} of (3.5) are $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{bs} \geq -\mu_{bs} \geq \dots \geq -\mu_2 \geq -\mu_1$. But the matrix T_{2s} is generated by $2s$ steps of the Block Lanczos method applied to \tilde{A} with initial matrix

$$x_1 = \begin{pmatrix} 0 \\ q_1 \end{pmatrix}.$$

If we let

$$G = (g_1, g_2, \dots, g_b),$$

then

$$G_1^t X_1 = \frac{1}{\sqrt{2}} V_1^t Q_1 .$$

The next theorem is a direct consequence of Theorem 2.1.

THEOREM 4.1. Assume that $\sigma_b > \sigma_{b+1}$. Let $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{bs}$ be the singular values of the $bs \times bs$ matrix J_s generated by the Block Bidiagonalization method with an initial $n \times b$ orthonormal matrix Q_1 . Suppose that the $b \times b$ matrix

$$W = \frac{1}{\sqrt{2}} V_1^t Q_1 ,$$

where

$$V_1 = (v_1, v_2, \dots, v_b),$$

is nonsingular so that its smallest singular value τ is positive.

Note that $\tau \leq 1/\sqrt{2}$. Then, for $i = 1, 2, \dots, b$, we have that

$$\sigma_i \geq \mu_i \geq \sigma_i - \epsilon_i^2 ,$$

where

$$\epsilon_i^2 = \frac{(\sigma_i - \sigma_1) \tan^2 \theta}{T_{2s-1}^2 \left(\frac{1 + r_i}{1 - r_i} \right)} ,$$

$$\theta = \cos^{-1} \tau$$

$$r_i = \frac{\sigma_i - \sigma_{b+1}}{\sigma_i + \sigma_1}$$

and

T_{2s-1} is the $(2s-1)$ -st Chebyshev polynomial of the first kind.

EXAMPLE 4.1.

Suppose that A is an $m \times n$ matrix with singular values $\sigma_1 = 1.0$, $\sigma_2 = 0.9$, $\sigma_3 = 0.5$, Let us apply the Block Bidiagonalization algorithm with $b = 2$ and $s = 5$. Then

$$r_1 = \frac{1.0 - 0.5}{1.0 + 1.0} = 0.25 ,$$

$$r_2 = \frac{0.9 - 0.5}{0.9 + 1.0} \doteq 0.21 ,$$

$$T_9 \left(\frac{1 + r_1}{1 - r_1} \right) \doteq T_9(1.67) \doteq 1.0 \times 10^4$$

and

$$T_9 \left(\frac{1 + r_2}{1 - r_2} \right) \doteq T_9(1.53) \doteq 3.7 \times 10^3 .$$

Let

$$\tau = \frac{1}{\sqrt{2}} \times 0.04 ,$$

so that

$$\tan^2 \theta \doteq 1249 .$$

Thus,

$$\epsilon_1^2 \doteq \frac{2.0 \times 1249}{1.0 \times 10^8} \doteq 2.5 \times 10^{-5}$$

and

$$\epsilon_2^2 \doteq \frac{1.9 \times 1249}{1.37 \times 10^7} \doteq 1.7 \times 10^{-4} .$$

The two greatest singular values μ_1 and μ_2 of J_s therefore satisfy the inequalities

$$\sigma_1 \geq \mu_1 \geq \sigma_1 - 2.5 \times 10^{-5}$$

and

$$\sigma_2 \geq \mu_2 \geq \sigma_2 - 1.7 \times 10^{-4}.$$

However, we suspect that the bounds of Theorem 4.1 may be gross overestimates. Suppose that

$$Q_1 = V_1.$$

Then we have

$$\tau = \frac{1}{\sqrt{2}}$$

and

$$\tan \theta = 1.$$

The last value is quite unsatisfactory for an initial matrix consisting of the correct singular vectors.

We seek to construct tighter bounds. From the matrix equations

$$A\bar{Q}_s = \bar{P}_s J_s \quad (3.12)$$

and

$$A^t \bar{P}_s = \bar{Q}_s J_s^t + \hat{Z}_{2s+1}, \quad (3.13)$$

we get

$$\begin{aligned} A^t A \bar{Q}_s &= A^t \bar{P}_s J_s \\ &= \bar{Q}_s J_s^t J_s + \hat{Z}_{2s+1} J_s. \end{aligned}$$

Thus,

$$A^t A \bar{Q}_s = \bar{Q}_s J_s^t J_s + (0, \dots, 0, Z_{2s+1} R_{2s}). \quad (4.1)$$

We observe that the matrix \bar{Q}_s is orthonormal and that the matrix $J_s^t J_s$ is block tridiagonal. It can be proved (cf. Chapter 2 and [14]) that the matrix equation (4.1) characterizes an application of the Block Lanczos method to the matrix $A^t A$ with the initial matrix Q_1 . Since the matrix $A^t A$ has eigenvalues $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_n^2$ and corresponding normalized eigenvectors $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$, we obtain the next result from Theorem 2.1.

THEOREM 4.2. Assume that $\sigma_b > \sigma_{b+1}$. Let $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{bs}$ be the singular values of the $bs \times bs$ matrix J_s generated by the Block Bidiagonalization method with an initial $n \times b$ orthonormal matrix Q_1 . Suppose that the $b \times b$ matrix

$$W = V_1^t Q_1,$$

where

$$V_1 = (\underline{v}_1, \underline{v}_2, \dots, \underline{v}_b),$$

is nonsingular so that its smallest singular value τ is positive.

Note that $\tau \leq 1$. Then, for $i = 1, 2, \dots, b$, we have that

$$\sigma_i \geq \mu_i \geq (\sigma_i^2 - \epsilon_i^2)^{1/2},$$

or

$$\sigma_i \geq \mu_i \geq \sigma_i - \frac{1}{2\sigma_i} \epsilon_i^2 \quad \text{for } \frac{\epsilon_i^2}{\sigma_i^2} < 1,$$

where

$$\epsilon_i^2 = \frac{(\sigma_i^2 - \sigma_n^2) \tan^2 \theta}{T_{s-1}^2 \left(\frac{1 + \gamma_i}{1 - \gamma_i} \right)},$$

$$\theta = \cos^{-1} \tau,$$

$$\gamma_i = \frac{\sigma_i^2 - \sigma_{b+1}^2}{\sigma_i^2 - \sigma_n^2}$$

and

T_{s-1} is the $(s-1)$ -st Chebyshev polynomial of the first kind.

EXAMPLE 4.2.

We use the same given data as in Example 4.1, with the additional assumption that $\sigma_n = 0.0$. Thus,

$$\sigma_1^2 = 1.00, \sigma_2^2 = 0.81, \sigma_3^2 = 0.25, \dots, \sigma_n^2 = 0.00,$$

so that

$$\gamma_1 = \frac{1.00 - 0.25}{1.00 - 0.00} = 0.75,$$

$$\gamma_2 = \frac{0.81 - 0.25}{0.81 - 0.00} \doteq 0.69,$$

$$T_4 \left(\frac{1 + \gamma_1}{1 - \gamma_1} \right) = T_4(7) \doteq 1.88 \times 10^4.$$

and

$$T_4 \left(\frac{1 + \gamma_2}{1 - \gamma_2} \right) \doteq T_4(5.45) \doteq 6.82 \times 10^3.$$

Also,

$$\tan^2 \theta = 624.$$

Consequently,

$$\epsilon_1^2 \doteq \frac{1.00 \times 624}{3.55 \times 10^8} \doteq 1.7 \times 10^{-6}$$

and

$$\epsilon_2^2 \doteq \frac{0.81 \times 624}{4.65 \times 10^7} \doteq 1.1 \times 10^{-5}.$$

As

$$\frac{1}{2\sigma_1} \epsilon_1^2 \doteq 9.0 \times 10^{-7}$$

and

$$\frac{1}{2\sigma_2} \epsilon_2^2 \doteq 6.8 \times 10^{-6},$$

the two greatest singular values μ_1 and μ_2 of J_s therefore satisfies the inequalities

$$\sigma_1 \geq \mu_1 \geq \sigma_1 - 9.0 \times 10^{-7}$$

and

$$\sigma_2 \geq \mu_2 \geq \sigma_2 - 6.8 \times 10^{-6}.$$

We observe that the bounds given by Theorem 4.2 are much smaller than those given by Theorem 4.1.

4.5. Iterating to Improve Accuracy

Let us restate our computational procedure. We use the Block Bidiagonalization method to generate a block bidiagonal matrix J_s of small order, and then apply standard techniques to compute the singular value decomposition of J_s . Our convergence test depends on the next two lemmas.

LEMMA 5.1 (Weinstein's Inequality) [10, p. 56].

Given a symmetric matrix B and a normalized vector \tilde{x} ,
if there is a scalar μ such that

$$\|B\tilde{x} - \mu\tilde{x}\| \leq \delta ,$$

then there is an eigenvalue λ of B within δ of μ .

LEMMA 5.2. [10, pp. 59-60].

Let B be an $\ell \times \ell$ symmetric matrix with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_\ell$ and let $\|B\tilde{x} - \lambda\tilde{x}\| \leq \delta$ for some normalized vector \tilde{x} . Suppose that

$$|\lambda - \lambda_i| \geq d > 0 \quad \text{for } i \neq j .$$

Then B has a normalized eigenvector \tilde{x}_j corresponding to λ_j such that

$$\|\tilde{x} - \tilde{x}_j\| \leq r(1 + r^2)^{1/2} ,$$

where $r = \delta/d$.

From equations (3.14) and (3.15), we obtain

$$\begin{pmatrix} 0 & A \\ A^t & 0 \end{pmatrix} \begin{pmatrix} p_i \\ q_i \end{pmatrix} = \mu_i \begin{pmatrix} p_i \\ q_i \end{pmatrix} + \begin{pmatrix} 0 \\ z_{2s+1}^h \tilde{i} \end{pmatrix} . \quad (5.1)$$

Thus, if

$$\|z_{2s+1}^h \tilde{i}\| \leq \delta \quad (5.2)$$

for some given tolerance δ , then there is a singular value σ_j of A such that

$$|\mu_i - \sigma_j| \leq \delta .$$

If it is also true that

$$|\mu_i - \sigma_k| \geq d > 0 \quad \text{for } k \neq j ,$$

where $\sigma_1, \sigma_2, \dots, \sigma_n$ are the singular values of A , then A has normalized left and right singular vectors \underline{u}_j and \underline{v}_j , respectively, corresponding to σ_j such that

$$\|\underline{p}_i - \underline{u}_j\|^2 + \|\underline{q}_i - \underline{v}_j\|^2 \leq \gamma^2(1 + \gamma^2) ,$$

where $\gamma = \delta/d$. For the more complicated case of multiple singular values, we refer the reader to an excellent paper of Stewart [13].

Suppose that our procedure has not computed all the k greatest singular values to the desired accuracy. In Theorem 4.2, the error bounds contain the term $\tan \theta$, where θ is the angle between the two subspaces spanned by the columns of V_1 and Q_1 . Let

$$\hat{Q}_1 = (\underline{q}_1, \underline{q}_2, \dots, \underline{q}_b) .$$

The results of the last section indicate that the matrix \hat{Q}_1 can be a better approximation to the matrix V_1 than the initial matrix Q_1 . Thus, we expect to compute more accurate approximations to the greatest singular values of A if we reapply the Block Bidiagonalization method with the initial matrix \hat{Q}_1 . This idea leads to the following iterative scheme, where we assume that b equals k .

ALGORITHM 5.1.

1. Let b , s and δ be given parameters and let Q_1 be a given $n \times b$ orthonormal matrix.
2. Repeat until all b greatest singular values have converged:
 - (a) Use the Block Bidiagonalization method with initial matrix Q_1 to generate the matrices J_s , \bar{P}_s and \bar{Q}_s .
 - (b) Compute the singular value μ_i and corresponding left and right singular vectors w_i and z_i , respectively, of J_s , for $i = 1, 2, \dots, bs$.
 - (c) Estimate the accuracy of μ_i as an approximation to σ_i , for $i = 1, 2, \dots, b$ (cf. inequality (5.2)).
 - (d) Let

$$Q_1 := (q_1, q_2, \dots, q_b),$$

where

$$q_i = \bar{Q}_s z_i \quad \text{for } i = 1, 2, \dots, b. \quad \square$$

This iterative algorithm provides a convenient means for estimating the accuracy of the computed singular values. We can show that the i -th column of the matrix Z_2 in the Block Bidiagonalization method computed in each iteration of Algorithm 5.1 after the first is the residual vector for the singular value μ_i computed in the previous iteration. It is thus possible to determine at the start of an iteration the accuracy of the singular values computed at the end of the previous iteration. We also observe that the matrices P_1 and R_2 of the Block Bidiagonalization method are readily available from the prior iteration.

However, once a few singular values and singular vectors have converged, we need not iterate with them any longer. Since the desired singular values may have different rates of convergence, as indicated by the error bounds of Theorems 4.1 and 4.2, we should modify Algorithm 5.1 so that (1) it does not iterate with those singular values and singular vectors that have converged, and (2) it allows the values of b and s to change from one iteration to the next. We are going to examine these issues in the next two sections.

4.6. Block Bidiagonalization Method With Reorthogonalization

Suppose that we are given accurate approximations to the k_0 greatest singular values $\sigma_1, \sigma_2, \dots, \sigma_{k_0}$ and corresponding singular vectors of the matrix A . Let P_0 and Q_0 be the $m \times k_0$ and $n \times k_0$ orthonormal matrices consisting of those given left and right singular vectors, respectively. We want now to compute the next $(k-k_0)$ greatest singular values of A .

Let us consider the matrix

$$\tilde{A} = \begin{pmatrix} 0 & A \\ A^t & 0 \end{pmatrix}. \quad (1.2)$$

If we apply Algorithm 2.1 to compute the eigenvalues $\pm \sigma_{k_0+1}$, $\pm \sigma_{k_0+2}$, \dots , $\pm \sigma_k$ of the matrix \tilde{A} , then we must maintain the orthogonality of the matrices X_1, X_2, \dots, X_s with respect to the eigenvectors corresponding to the eigenvalues $\pm \sigma_1, \pm \sigma_2, \dots, \pm \sigma_{k_0}$ (cf. [14]). But these eigenvectors are accurately approximated by the columns of the matrix (see [8, Chap. 3])

$$\frac{1}{\sqrt{2}} \begin{pmatrix} P_0 & P_0 \\ Q_0 & -Q_0 \end{pmatrix}.$$

We have that

$$I - \frac{1}{\sqrt{2}} \begin{pmatrix} P_0 & P_0 \\ Q_0 & -Q_0 \end{pmatrix} \frac{1}{\sqrt{2}} \begin{pmatrix} P_0 & P_0 \\ Q_0 & -Q_0 \end{pmatrix}^t = \begin{pmatrix} I - P_0 P_0^t & 0 \\ 0 & I - Q_0 Q_0^t \end{pmatrix}.$$

Therefore, if we want to apply the Block Bidiagonalization method to compute the $(k-k_0)$ next greatest singular values of A , we need to maintain the orthogonality of the matrices P_1, P_2, \dots, P_s with respect to P_0 , and the orthogonality of the matrices Q_1, Q_2, \dots, Q_s with respect to Q_0 . We are thus computing the $(k-k_0)$ greatest

singular values of the "deflated" matrix

$$\hat{A} = (I - P_0 P_0^t) A (I - Q_0 Q_0^t) \quad (6.1)$$

Suppose that

$$\|A Q_0 - P_0 \Sigma_0\|_F^2 + \|A^t P_0 - Q_0 \Sigma_0\|_F^2 \leq \delta^2, \quad (6.2)$$

where $\|\cdot\|_F$ is the Frobenius matrix norm and $\Sigma_0 = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_{k_0})$.

Note that if both matrices P_0 and Q_0 had been computed by the Block Lanczos procedure (cf. equation (3.12)), then

$$\|A Q_0 - P_0 \Sigma_0\| = m k_0 O(\epsilon), \quad (6.3)$$

where ϵ is the machine precision. Using arguments similar to those in Underwood's thesis [14, pp. 62-66], we get from (6.2) that the $(k-k_0)$ greatest singular values of \hat{A} differ from the $(k-k_0)$ next greatest singular values of A by quantities which are less than $|\delta|$ in modulus.

Now, the Block Bidiagonalization method may be numerically unstable. Although the matrices P_i 's and Q_i 's form two sequences

of orthonormal matrices in exact arithmetic, they lose orthogonality rapidly in practice due to the loss of figures in the computation of the matrices Z_i 's and W_i 's. To maintain stability in the scheme, we choose to reorthogonalize the matrices P_i and Q_i with respect to all the previous P_j 's and Q_j 's, respectively (cf. [14]).

We have incorporated the two reorthogonalization procedures into the next algorithm.

ALGORITHM 6.1. (Block Bidiagonalization method with Reorthogonalization)

1. Let Q_1 be a given $n \times b$ orthonormal matrix.

Compute

$$W_1 := A Q_1 .$$

Orthogonalize W_1 with respect to P_0 . Factorize W_1 so that

$$W_1 := P_1 R_2 ,$$

where P_1 is orthonormal and R_2 is upper triangular.

2. For $i = 2, 3, \dots, s$, do

(a) Compute

$$Z_i := A^t P_{i-1} - Q_{i-1} R_{2i-2}^t .$$

Orthogonalize Z_i with respect to Q_0, Q_1, \dots, Q_{i-1} .

Factorize Z_i so that

$$Z_i := Q_i R_{2i-1} ,$$

where Q_i is orthonormal and R_{2i-1} is upper triangular.

(b) Compute

$$W_i := AQ_i - P_{i-1} R_{2i-1}^t .$$

Orthogonalize W_i with respect to P_0, P_1, \dots, P_{i-1} .

Factorize W_i so that

$$W_i := P_i R_{2i} ,$$

where P_i is orthonormal and R_{2i} is upper triangular. \square

In case of rank deficiency of the matrix Z_i or W_i , we apply a remedial procedure similar to that for Algorithm 3.1.

4.7. Iterative Block Lanczos Method

We should point out that the reorthogonalization process of the last section not only requires a large number of arithmetic operations but also requires that each of the P_j 's and Q_j 's be in memory during each step of Algorithm 6.1. Since m and n are large numbers in this application, the available computer memory places an upper bound c on the product bs . It is then necessary to determine optimal values for b and s subject to this upper bound constraint. But the error bounds of Theorem 4.2 indicate that we need accurate knowledge of the singular value spectrum of the given matrix, which is precisely the same information we are trying to obtain.

A good initial choice of the block size b is the number k of singular values to be computed (see [2] and [14]). This may not be

the best choice, as we can see from Example 4 of the next section. Our experiments have shown that it seldom pays to have $b > k$ (cf. Example 3 of the next section). Underwood [14] made the same observation, but he used a different way to update the block size.

Having chosen b , we compute the number s of blocks by

$$s := \left\lfloor \frac{c}{b} \right\rfloor ,$$

where $\lfloor \alpha \rfloor$ denotes the integer part of a real number α . If s is less than 2, we compute

$$b := \left\lfloor \frac{c}{2} \right\rfloor$$

and

$$s := \left\lfloor \frac{c}{b} \right\rfloor .$$

The last computation is necessary so that s would equal 3 if the value of c were 3.

Let us describe how we update the values of c , b , s and k . Suppose that k_0 singular values and associated singular vectors have converged in an iteration. Then

$$c := c - k_0$$

because those computed singular vectors must reside in the computer memory for the reorthogonalization process. Now, if $b \geq k$, then we decrease the value of b by k_0 , i.e.

$$b := b - k_0 ;$$

otherwise we choose the new block size as the smaller value of the old block size and the number of singular values left to be computed, i.e.

$$b := \min(b, k - k_0) .$$

We update the value of k by

$$k := k - k_0 .$$

The new value for s is then computed in the same manner as described in the previous paragraph, with the same modification to the value of b if necessary. Our scheme for updating b differs from Underwood's [14] only in the case where $b > k$.

ALGORITHM 7.1. (Iterative Block Lanczos method)

1. Let c, b, s and δ be given parameters and let Q_1 be a given $n \times b$ orthonormal matrix. The matrices P_0 and Q_0 are null.
2. Repeat until all k singular values have converged:
 - (a) Use Algorithm 6.1 with initial matrix Q_1 to generate the matrices J_s, \bar{P}_s and \bar{Q}_s .
 - (b) Compute the singular value μ_i and corresponding left and right singular vectors w_i and z_i , respectively, of \bar{J}_s , for $i = 1, 2, \dots, bs$.

- (c) Estimate the accuracy of μ_i as an approximate singular value, for $i = 1, 2, \dots, bs$. Assume that k_0 singular values have converged.
- (d) Update the values of c , b and s .
- (e) For $i = 1, 2, \dots, k_0 + b$, compute

$$p_i := \bar{P}_s w_i$$

and

$$g_i := \bar{Q}_s z_i .$$

- (f) Let

$$P_0 := (P_0 | p_1, \dots, p_{k_0})$$

and

$$Q_0 := (Q_0 | g_1, \dots, g_{k_0}) .$$

- (g) Let

$$Q_1 := (g_{k_0+1}, g_{k_0+2}, \dots, g_{k_0+b}) . \quad \square$$

4.8. Test Examples

Rectangular diagonal matrices are chosen for all our examples. Such matrices are sufficiently general for the Lanczos method which does not transform the given matrix. We can thus specify the singular value spectrum and study the behavior of the algorithm as a function of b and s .

A set of FORTRAN routines has been written to implement Algorithm 7.1 (see [5]). We ran our tests on an IBM 370/168 computer

at the Stanford Linear Accelerator Center. The code was compiled by the H EXTENDED compiler with optimization level 2.

The computed singular value μ_i and associated singular vectors p_i and q_i are accepted if they satisfy the inequality

$$(\|Aq_i - \mu_i p_i\|^2 + \|A^t p_i - \mu_i q_i\|^2)^{1/2} \leq 10^{-3}.$$

In Sections 4.5 and 4.6, we have described a way to test for convergence with very little additional work. We have chosen the upper bound c for the product bs to be 12.

The following notations are used in the examples:

$\sigma_1, \sigma_2, \sigma_3, \dots$ are the computed singular values in the order of convergence.

Iter = total number of iterations.

Time = machine execution time in seconds.

$m' - n = m \times 10^{-n}$.

EXAMPLE 1.

A is a 905×904 matrix with diagonal elements -1.00, -0.99, -0.98, and 0.000, 0.001, \dots , 0.900.

k = 3	b = 3
σ_1	1.00 - 4' - 7
σ_2	0.99 - 9' - 9
σ_3	0.98 - 3' - 6
Iter	5
Time	11.36

We note that the computed singular values, as Rayleigh quotients, are accurate to twice the number of digits of the error tolerance.

EXAMPLE 2.

A is a 905×904 matrix with diagonal elements $-1.000, -0.999, -0.998$, and $0.9000, 0.9001, \dots, 0.9900$. This example is essentially the same as the previous one except that the gaps between the singular values have been reduced by a factor of 10.

$k = 3$	$b = 3$
σ_1	$1.000 - 9' - 7$
σ_2	$0.999 - 4' - 8$
σ_3	$0.998 - 2' - 8$
Iter	6
Time	13.83

The first two examples illustrate the fact that the convergence rate of the Lanczos algorithm depends on the relative spread of the singular values (cf. Theorems 4.1 and 4.2).

EXAMPLE 3.

A is an 806×805 matrix with diagonal elements
1.0, -1.0, 0.9, -0.9, and 0.000, 0.001, ... , 0.800.

k = 1	b = 1	b = 2
σ_1	1.0 - 3' - 10	1.0 - 1' - 12
Iter	1	3
Time	2.35	6.52

k = 2	b = 1	b = 2
σ_1	1.0 - 3' - 10	1.0 - 1' - 12
σ_2	0.9 - 3' - 16	1.0 - 1' - 11
Iter	2	3
Time	4.27	6.52

EXAMPLE 3 continued

$k = 3$	$b = 1$	$b = 2$	$b = 3$	$b = 4$
σ_1	$1.0 - 3' - 10$	$1.0 - 1' - 12$	$1.0 - 3' - 8$	$1.0 - 1' - 7$
σ_2	$0.9 + 3' - 16$	$1.0 - 1' - 11$	$1.0 - 4' - 9$	$1.0 - 8' - 7$
σ_3	$1.0 + 0' - 16$	$0.9 - 1' - 9$	$0.9 - 3' - 9$	$0.9 - 6' - 9$
Iter	5	4	4	5
Time	9.60	8.28	8.23	10.53

$k = 4$	$b = 1$	$b = 2$	$b = 3$	$b = 4$
σ_1	$1.0 - 3' - 10$	$1.0 - 1' - 12$	$1.0 - 3' - 8$	$1.0 - 1' - 7$
σ_2	$0.9 + 3' - 16$	$1.0 - 1' - 11$	$1.0 - 6' - 8$	$1.0 - 8' - 7$
σ_3	$1.0 + 0' - 16$	$0.9 - 3' - 8$	$0.9 - 2' - 7$	$0.9 - 6' - 9$
σ_4	$0.9 + 6' - 14$	$0.9 - 7' - 8$	$0.9 + 8' - 10$	$0.9 - 2' - 7$
Iter	7	5	5	5
Time	12.70	10.00	9.59	10.38

These tests illustrate a couple of important points. First, the Lanczos method does not always compute the greatest singular values (cf. the case where $k = 2$ and $b = 1$). Second, the supposition of Theorems 4.1 and 4.2 that $\sigma_b > \sigma_{b+1}$ is not necessary for the convergence of the Lanczos method.

EXAMPLE 4.

A is a 902×901 matrix with diagonal elements 0.000, 0.001, ..., 0.900.

$k = 3$	$b = 1$	$b = 2$	$b = 3$
σ_1	0.900 - 5' - 6	0.900 - 1' - 5	0.900 - 9' - 6
σ_2	0.899 + 3' - 6	0.899 - 2' - 6	0.899 - 9' - 6
σ_3	0.898 - 2' - 5	0.898 - 4' - 7	0.898 - 7' - 5
Iter	13	27	23
Time	28.73	61.85	52.53

This example of a dense singular value spectrum is one in which the point algorithm ($b = 1$) works better than the block algorithm ($b \geq 2$).

Our Iterative Block Lanczos method is therefore a good procedure for computing a few greatest singular values of a matrix. A block method with an appropriate block size can (1) have a fast convergence rate, and (2) handle well the case of multiple singular values (see Example 3). For problems where the given matrix has to be read from secondary storage, economics may dictate that we multiply the matrix into a block of vectors and thus choose a block method.

REFERENCES

- [1] Cullum, J., "The simultaneous computation of a few algebraically largest and smallest eigenvalues of a large, sparse, symmetric matrix," Report RC 6827, IBM Thomas J. Watson Research Center, Yorktown Heights (1977).
- [2] Cullum, J., and Donath, W. E., "A block Lanczos algorithm for computing the q algebraically largest eigenvalues and a corresponding eigenspace of large, sparse, real symmetric matrices," Proc. 1974 IEEE Conf. on Decision and Control, Phoenix, Arizona (1974), 505-509.
- [3] Garbow, B.S., Boyle, J.M., Dongarra, J.J., and Moler, C.B., Matrix Eigensystem Routines--EISPACK Guide Extension, Springer-Verlag, Berlin (1977).
- [4] Golub, G.H., and Luk, F.T., "Singular value decomposition: applications and computations," ARO Report 77-1, Transactions of the 22nd Conference of Army Mathematicians (1977), 577-605.
- [5] Golub, G.H., Luk, F.T., and Overton, M.L., "A block Lanczos method to compute the singular values and corresponding singular vectors of a matrix," Report STAN-CS-77-635, Computer Science Dept., Stanford University (1977).
- [6] Golub, G.H., and Reinsch, C., "Singular value decomposition and least squares solutions," Numer. Math. 14 (1970), 403-420.
- [7] Golub, G.H., and Underwood, R., "The Block Lanczos method for computing eigenvalues," in Mathematical Software III, (J. R. Rice, Ed.), Academic Press, New York (1977), 361-377.
- [8] Lanczos, C., Linear Differential Operators, Van Nostrand, London (1961).
- [9] Lewis, J.G., "Algorithms for sparse matrix eigenvalue problems," Ph.D. thesis. Report STAN-CS-77-595, Computer Science Dept., Stanford University (1977).
- [10] Ortega, J.M., Numerical Analysis: A Second Course, Academic Press, New York (1972).
- [11] Ruhe, A., "Implementation aspects of Band Lanczos algorithms for computation of eigenvalues of large sparse symmetric matrices," Report, Dept. of Mathematics, University of California, San Deigo (February 1978).
- [12] Smith, B.T., Boyle, J.M., Dongarra, J.J., Garbow, B.S., Ikebe, Y., Klema, V.C., and Moler, C.B., Matrix Eigensystem Routines--EISPACK Guide, Second Edition, Springer-Verlag, Berlin (1976).

- [13] Stewart, G.W., "Error and perturbation bounds for subspaces associated with certain eigenvalue problems," SIAM Review 15 (1973), 727-764.
- [14] Underwood, R., "An iterative Block Lanczos method for the solution of large sparse symmetric eigenproblems," Ph.D. thesis, Report STAN-CS-75-496, Computer Science Dept., Stanford University (1975).
- [15] Wilkinson, J.H., and Reinsch, C., Handbook for Automatic Computation, Volume II, Linear Algebra, Part 2, Springer-Verlag, Berlin (1971).

V. COMPUTING THE SINGULAR VALUE DECOMPOSITION ON THE ILLIAC IV

5.1. Introduction

We study the computation of the singular value decomposition on the ILLIAC IV computer. Suppose that we have a real $m \times n$ matrix A . Its singular value decomposition (SVD) can be defined as

$$A = U \Sigma V^t \quad (1.1)$$

with

$$U^t U = V^t V = I_k \quad \text{and} \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_k),$$

where

$$k = \min(m, n). \quad (1.2)$$

The matrices U and V consist of the orthonormalized eigenvectors associated with the k largest eigenvalues of AA^t and $A^t A$, respectively. The diagonal elements of Σ are the non-negative square roots of the k largest eigenvalues of AA^t , and are called the singular values. We assume that

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0 \quad \text{and} \quad \sigma_{r+1} = \dots = \sigma_k = 0, \quad (1.3)$$

i.e. the rank of A equals r . An alternative definition of the singular value decomposition is

$$A = U_r \Sigma_r V_r^t \quad (1.4)$$

with

$$U_r^t U_r = V_r^t V_r = I_r \quad \text{and} \quad \Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r).$$

The singular value decomposition is a very useful matrix decomposition (see [7]). Various methods have been proposed for its computation. The standard method was introduced by Golub and Kahan in 1965 [6]. They use first the Householder transformation to bidiagonalize the given matrix, and then the QR method to compute the singular values of the resultant bidiagonal form. Their method superceded a one-sided orthogonalization method given by Hestenes in 1958 [10]. Hestenes' method is, however, easily adaptable to special purpose computations. It was suggested by Chartres (1962) [2] for a computer with a magnetic backing store, and was implemented on a mini-computer by Nash [12]. In this chapter, we study the implementation of Hestenes' method on the ILLIAC IV computer and show that the method makes very efficient use of the parallel computing abilities of the ILLIAC machine.

We are going to use the Frobenius norm for matrices, i.e.

$$\|A\| = \|A\|_F = \left(\sum_{i,j} a_{ij}^2 \right)^{1/2} \quad \text{for } A \equiv (a_{ij}),$$

and the Euclidean norm for vectors, i.e.

$$\|\underline{x}\| = \|\underline{x}\|_2 = (\underline{x}^t \underline{x})^{1/2}.$$

5.2. The ILLIAC IV Computer

The ILLIAC IV computer was built by the Burroughs Corporation and is located at NASA/Ames Research Center, Moffett Field, California. The computer consists of 64 synchronous processing elements (PE's)

under the direction of a single control unit (CU). Each PE has 2048 words of 64-bit memory with an access time of 188 nanoseconds, and is capable of performing a general floating-point operation in about 1.7 microseconds and a typical bookkeeping operation in about 1.2 microseconds. The PE instruction set is similar to that of conventional machines with two exceptions. First, each PE can communicate data to four other PE's through routing instructions. Second, the PE's can set their own mode registers to effectively disable or enable themselves. The CU takes about 0.7 microseconds to perform a bookkeeping operation.

The main memory of the ILLIAC is logically a 16-million word drum, which is divided into 52 bands and has a 40 millisecond rotation period. Data transfers to or from the PE memory are program initiated and are performed in blocks of 1024 words. The transfer time for 1024 words is about 66 microseconds; it takes about 4.2 milliseconds to refresh half the PE memory.

A floating-point number on the ILLIAC consists of a 1-bit sign, a 15-bit exponent to the radix 2, and a normalized 48-bit mantissa. The machine precision ϵ is thus about 3.55×10^{-15} . A fixed-point number has a 1-bit sign and a 48-bit mantissa.

5.3. Programming Languages for the ILLIAC

There are three languages available for programming the ILLIAC; its assembly language, ASK; a FORTRAN-like language, CFD [15]; and an ALGOL 60-like language, GLYPNIR [11]. Both CFD and

GLYPNIR do not hide the basic 64-wide architecture of the ILLIAC. We must restructure our data and algorithm so that the computation can be done in parallel in "strips" of width 64 or less.

Let us briefly describe the data declarations in GLYPNIR. The PE memory of the ILLIAC can be viewed as a two-dimensional structure where each word can be addressed by an ordered pair which specifies the PE memory module and the address within that module. A group of 64 words, each in a different module but each having the same address within its module, is called a superword or sword. We can divide the variable types in GLYPNIR into two major categories. The first represents words or vectors of words; they are called the CU variables. The second represents swords or vectors of swords; they are called the PE variables. There are also the Boolean variables and the so-called ADB variables.

A sword vector of length n represents an indexable vector of swords. It is thus in some sense an $n \times 64$ array. A GLYPNIR program cannot directly handle two-dimensional arrays whose row and column dimensions exceed 64.

5.4. A Row Orthogonalization Method

There are two reasons why the standard SVD method of Golub and others (see [6] and [8]) may be undesirable on a parallel processor. First, although the Householder transformation is inherently parallel, the effective vector length decreases at each step, causing inefficiencies.

Second, the parallel QR method [14] may be numerically unstable (see [9]). In contrast, the one-sided orthogonalization method of Hestenes [10] is easily adaptable to computation on a parallel machine.

The method of Hestenes consists of generating an orthogonal matrix V such that the non-null column vectors of the matrix

$$H = AV$$

are mutually orthogonal and non-increasing in norm. The nonzero columns of H are then normalized so that

$$H = (U_r | 0) \left(\begin{array}{c|c} \Sigma_r & 0 \\ \hline 0 & 0 \end{array} \right)$$

with

$$U_r^t U_r = I_r \quad \text{and} \quad \Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r).$$

Consequently,

$$A = U_r \Sigma_r V_r^t, \quad (1.4)$$

where

V_r is an $n \times r$ matrix consisting of the first r columns of V .

Nash [12] followed Hestenes' approach, but Chartres [2] chose to orthogonalize the rows of the given matrix A . We have decided on the row orthogonalization scheme, for it is easily adaptable to solving overdetermined linear equations.

We aim to generate an orthogonal matrix U^t so that the non-null row vectors of the matrix

$$K = U^t A$$

are mutually orthogonal and non-increasing in norm. We then normalize the nonzero rows of K to obtain

$$K = \left(\begin{array}{c|c} \Sigma_r & 0 \\ \hline 0 & 0 \end{array} \right) \left(\begin{array}{c} V_r^t \\ 0 \end{array} \right)$$

with

$$\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r) \quad \text{and} \quad V_r^t V_r = I_r.$$

It follows that

$$A = U \Sigma_r V_r^t, \quad (1.4)$$

where

U_r is an $m \times r$ matrix consisting of the first r columns of U .

We are going to construct the matrix U as a product of plane rotations. Let us write the matrix A as

$$A = \begin{pmatrix} a_1^t \\ a_2^t \\ \vdots \\ a_m^t \end{pmatrix}, \quad (4.1)$$

where

a_i^t is an $1 \times n$ row vector, for $i = 1, 2, \dots, m$.

Given any two rows \tilde{a}_i^t and \tilde{a}_j^t , with $i < j$. We would consider them orthogonal if

$$\|\tilde{a}_i\| < \epsilon \quad \text{or} \quad \|\tilde{a}_j\| < \epsilon, \quad (4.2)$$

where ϵ is the machine precision, or if

$$\frac{\tilde{a}_i^t \tilde{a}_j^t}{\|\tilde{a}_i\| \|\tilde{a}_j\|} < \tau, \quad (4.3)$$

where τ is a previously chosen tolerance. We do not transform orthogonal rows, but would permute them if

$$\|\tilde{a}_i\| < \|\tilde{a}_j\|.$$

Suppose now that the two given rows do not satisfy the orthogonality condition (4.2) or (4.3). Let us consider the action of a plane rotation:

$$\begin{pmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{pmatrix} \begin{pmatrix} \tilde{a}_i^t \\ \tilde{a}_j^t \end{pmatrix} = \begin{pmatrix} \hat{a}_i^t \\ \hat{a}_j^t \end{pmatrix}. \quad (4.4)$$

The idea is to choose φ such that

$$\hat{a}_i^t \hat{a}_j^t = 0 \quad \text{and} \quad \|\hat{a}_i\| \geq \|\hat{a}_j\|.$$

The second condition ensures that the computation always proceeds towards an ordering of row norms. Following Nash [12], we let

$$\begin{aligned}\alpha &= 2\tilde{a}_i^t \tilde{a}_j \\ \beta &= \|\tilde{a}_i\|^2 - \|\tilde{a}_j\|^2,\end{aligned}\tag{4.5}$$

and

$$\gamma = (\alpha^2 + \beta^2)^{1/2}.$$

Note that γ is positive since α is nonzero. Then, if β is positive, we compute

$$\cos \varphi = \left(\frac{\gamma + \beta}{2\gamma} \right)^{1/2} \quad \text{and} \quad \sin \varphi = \frac{\alpha}{2\gamma \cos \varphi}, \tag{4.6}$$

otherwise we let

$$\sin \varphi = \left(\frac{\gamma - \beta}{2\gamma} \right)^{1/2} \quad \text{and} \quad \cos \varphi = \frac{\alpha}{2\gamma \sin \varphi}. \tag{4.7}$$

In (4.4), we could use the Fast Givens transformations (see [4]) which requires only $2n$ multiplications, an apparent 50% work reduction. But a heavy overhead in maintaining the scaling factors eats up the savings unless the row length n is moderately large [16].

As in the traditional Jacobi algorithm, the plane rotations are performed in a set sequence called a sweep, which consists of the $[m(m-1)]/2$ plane rotations on the row pairs

$$(1,2), (1,3), \dots, (1,m), (2,3), \dots, (2,m), (3,4), \dots, (m-1,m).$$

The iterative procedure terminates if (1) all the rows are pairwise orthogonal, and (2) no row permutations have occurred, in one complete sweep.

Our orthogonalization method is in essence the Jacobi method implicitly applied to the matrix AA^t to compute its eigenvalues. We can refer to the literature [17] for the convergence properties of our method. We see that the convergence is quadratic and takes the order of 6 to 10 sweeps, i.e. from $3m^2$ to $5m^2$ plane rotations (see [13]).

We now present our method in its entirety. Two Boolean variables are introduced:

withu : true if matrix U is desired, false otherwise.

withv : true if matrix V is desired, false otherwise.

We make the arbitrary choice that

$$\tau = 10^{-12},$$

and write the matrices U and V as

$$U = (\underline{u}_1, \underline{u}_2, \dots, \underline{u}_m), \quad (4.8)$$

and

$$V = (\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n).$$

ALGORITHM 1 (SVD).

I. Initialize:

(1) Let

$$\epsilon := 3.55 \times 10^{-15},$$

$$\tau := 10^{-12},$$

and

$$c := 0.$$

(2) If (withu) then let $U := I$.

II. Repeat until $c = [m(m-1)]/2$:

(1) Let $c := 0$

(2) For $(i, j) := (1, 2), (1, 3), \dots, (1, m), (2, 3), \dots, (2, m),$
 $(3, 4), \dots, (3, m), \dots, (m-1, m)$ do

If $\|a_j\| < \epsilon$ then

(a) Let $c := c + 1$.

Else if $\|a_i\| < \epsilon$ then

(a) Exchange a_i and a_j .

(b) If (withu) then exchange u_i and u_j .

Else if

$$\frac{a_i^t a_j}{\|a_i\| \|a_j\|} < \tau$$

then

(a) Let $c := c + 1$.

(b) If $\|a_i\| < \|a_j\|$ then

(i) Exchange a_i and a_j .

(ii) If (withu) then exchange u_i and u_j .

Else

(a) Compute

$$\alpha := 2a_i^t a_j$$

and

$$\beta := \|a_i\|^2 - \|a_j\|^2,$$

$$\gamma := (\alpha^2 + \beta^2)^{1/2}.$$

(b) If $\beta > 0$ then

(i) Compute

$$\cos \varphi := \left(\frac{r + \beta}{2r} \right)^{1/2},$$

and

$$\sin \varphi := \frac{\alpha}{2r \cos \varphi}.$$

Else

(i) Compute

$$\sin \varphi := \left(\frac{r - \beta}{2r} \right)^{1/2},$$

and

$$\cos \varphi := \frac{\alpha}{2r \sin \varphi}.$$

(c) Compute

$$\tilde{w} := \cos \varphi \cdot \tilde{a}_i + \sin \varphi \cdot \tilde{a}_j,$$

$$\tilde{a}_j := -\sin \varphi \cdot \tilde{a}_i + \cos \varphi \cdot \tilde{a}_j,$$

and

$$\tilde{a}_i := \tilde{w}.$$

(d) If (withu) then

$$\tilde{z} := \cos \varphi \cdot \tilde{u}_i + \sin \varphi \cdot \tilde{u}_j,$$

$$\tilde{u}_j := -\sin \varphi \cdot \tilde{u}_i + \cos \varphi \cdot \tilde{u}_j,$$

and

$$\tilde{u}_i := \tilde{z}.$$

III. Compute singular values:

(1) Let $i := 1$.

(2) Repeat until $i > m$ or $\|\tilde{a}_i\| < \epsilon$:

(a) Let $\sigma_i := \|\tilde{a}_i\|$.

(b) Compute $\tilde{v}_i := \frac{\tilde{a}_i}{\|\tilde{a}_i\|}$.

(c) Let $i := i + 1$.

(3) Let $r := i - 1$.

We wish to compare the required work of Algorithm 1 and the Golub-Reinsch method [8]. One sweep of Algorithm 1 takes about $(7n + 4m)m^2/2$ multiplications if the U matrix is desired, and $7nm^2/2$ multiplications otherwise. We assume that

$$m \leq n ,$$

for we can compute the SVD of A^t if $m > n$. We further suppose that our Jacobi-like method takes 8 sweeps to converge and that only two QR steps are required per singular value for the Golub-Reinsch algorithm (cf. [1]). The following table gives the number of multiplications required by the two methods in four different cases.

Matrices Desired	Algorithm 1	Golub-Reinsch
U_r, Σ_r, V_r	$28m^2n + 16m^3$	$7m^2n + 11m^3/3$
U_r, Σ_r	$28m^2n + 16m^3$	$7m^2n - m^3$
Σ_r, V_r	$28m^2n$	$2m^2n + 4m^3$
Σ_r	$28m^2n$	$2m^2n - 2m^3/3$

We see that Algorithm 1 is about four times slower than the standard SVD algorithm in computing the full singular value decomposition. However, the special architecture of the ILLIAC IV computer can reduce the number of required multiplications by an asymptotic factor of 64. Our Jacobi-like algorithm is therefore very efficient on a parallel computer. We should mention that Chan [1] described a modified Golub-Reinsch algorithm that could save up to 50% of machine execution time if $m \ll n$.

5.5. Least Squares Solutions

We refer to an $n \times m$ matrix S as the pseudoinverse of an $m \times n$ matrix A if S satisfies the following four conditions:

$$\begin{aligned} \text{(i)} \quad & ASA = A, \\ \text{(ii)} \quad & SAS = S, \\ \text{(iii)} \quad & (AS)^t = AS, \\ \text{and (iv)} \quad & (SA)^t = SA. \end{aligned} \tag{5.1}$$

Let us denote the pseudoinverse by A^+ . Now, we have defined the singular value decomposition of A as

$$A = U \Sigma V^t \tag{1.1}$$

We can therefore write the pseudoinverse as

$$A^+ = V \Omega U^t, \tag{5.2}$$

where Ω is an $n \times m$ matrix to be determined. Using the four conditions of (5.1), we easily determine that Ω is given uniquely by

$$\Omega = \left(\begin{array}{ccc|ccc} \sigma_1^{-1} & & & & & \\ & \cdot & & \bigcirc & & \\ & & \cdot & & & \bigcirc \\ & \bigcirc & & & \sigma_r^{-1} & \\ \hline & & & \bigcirc & & \bigcirc \end{array} \right). \tag{5.3}$$

An important application of the pseudoinverse is in solving linear equations

$$AX = B, \quad (5.4)$$

where B is a given $m \times s$ matrix. We assume that

$$m \geq n, \quad (5.5)$$

and we seek an $n \times s$ matrix X such that

$$\|B - AX\| = \min. \quad (5.6)$$

The solution matrix X is not unique unless the matrix A is of full rank. We therefore impose the condition that we want the matrix \hat{X} of minimum norm in the solution space. It is well known (see, e.g., [6]) that \hat{X} is unique and is given by

$$\hat{X} = A^+ B. \quad (5.7)$$

Thus, we have that

$$\hat{X} = V\Omega C, \quad (5.8)$$

where

$$C = U^t B.$$

The matrix C can be generated by applying to the rows of B those plane rotations that we use to orthogonalize the rows of A . It is unnecessary to accumulate the plane rotations.

We now present an algorithm based on Algorithm SVD for computing the minimum norm solution to the overdetermined system (5.4). There is an input parameter "cutoff." Our method sets to zero all those singular values of A that are smaller than "cutoff." The $1 \times s$ row vector \tilde{b}_i^t denotes the i -th row of B , for $i = 1, 2, \dots, m$.

ALGORITHM 2. (MINFIT)

I. Initialize:

(1) Let

$$\epsilon := 3.55 \times 10^{-15},$$

$$\tau := 10^{-12}$$

and

$$c := 0.$$

II. Repeat until $c = [m(m-1)]/2$:

(1) Let $c := 0$.

(2) For $(i, j) := (1, 2), (1, 3), \dots, (1, m), (2, 3), \dots, (2, m),$
 $(3, 4), \dots, (3, m), \dots, (m-1, m)$ do

If $\|\tilde{a}_j\| < \epsilon$ then

(a) Let $c := c + 1$.

Else if $\|\tilde{a}_i\| < \epsilon$ then

(a) Exchange \tilde{a}_i and \tilde{a}_j .

(b) Exchange \tilde{b}_i and \tilde{b}_j .

Else if

$$\frac{\tilde{a}_i^t \tilde{a}_j}{\|\tilde{a}_i\| \|\tilde{a}_j\|} < \tau$$

then

(a) Let $c := c + 1$.

(b) If $\|\tilde{a}_i\| < \|\tilde{a}_j\|$ then

(i) Exchange \tilde{a}_i and \tilde{a}_j .

(ii) Exchange \tilde{b}_i and \tilde{b}_j .

Else

(a) Compute

$$\alpha := 2\tilde{a}_i^t \tilde{a}_j$$

$$\beta := \|\tilde{a}_i\|^2 - \|\tilde{a}_j\|^2,$$

and

$$\gamma := (\alpha^2 + \beta^2)^{1/2}.$$

(b) If $\beta > 0$ then

(i) Compute

$$\cos \varphi := \left(\frac{\gamma + \beta}{2\gamma} \right)^{1/2}$$

and

$$\sin \varphi := \frac{\alpha}{2\gamma \cos \varphi} .$$

Else

(i) Compute

$$\sin \varphi := \left(\frac{\gamma - \beta}{2\gamma} \right)^{1/2}$$

and

$$\cos \varphi := \frac{\alpha}{2\gamma \sin \varphi} .$$

(c) Compute

$$\tilde{w} := \cos \varphi \cdot \tilde{a}_i + \sin \varphi \cdot \tilde{a}_j ,$$

$$\tilde{a}_j := -\sin \varphi \cdot \tilde{a}_i + \cos \varphi \cdot \tilde{a}_j ,$$

$$\text{and } \tilde{a}_i := \tilde{w} .$$

(d) Compute

$$\tilde{z} := \cos \varphi \cdot \tilde{b}_i + \sin \varphi \cdot \tilde{b}_j ,$$

$$\tilde{b}_j := -\sin \varphi \cdot \tilde{b}_i + \cos \varphi \cdot \tilde{b}_j ,$$

and

$$\tilde{b}_i := \tilde{z} .$$

III. Compute least squares solution:

(1) Let

$$V := 0 ,$$

$$Y := 0 ,$$

and

$$i := 1 .$$

(2) Repeat until $i > n$ or $\|\tilde{a}_i\| < \text{cutoff}$:

(a) Compute

$$\tilde{v}_i := \frac{\tilde{a}_i}{\|\tilde{a}_i\|}$$

and

$$\tilde{b}_i := \frac{\tilde{b}_i}{\|\tilde{a}_i\|}.$$

(b) Let $i := i + 1$.

(3) Let $r := i - 1$.

(4) Compute $X := VB$.

One sweep of Algorithm 2 takes about $7m^2n/2$ multiplications. We neglect the terms involving s because $s \ll n$ in most applications. With the same assumptions on convergence rates as in the last section, the required work for our MINFIT algorithm is about $28m^2n$ multiplications while that for a similar method based on the Golub-Reinsch algorithm (see [8]) is about $2mn^2 + 4n^3$ multiplications. If $m \gg n$, it saves work to first reduce the regression matrix A to upper triangular form using Householder transformations, before applying the MINFIT algorithm (cf. [1]). Such a two-stage scheme requires about $mn^2 + 83n^2/3$ multiplications. Fortunately, the parallel computing abilities of the ILLIAC machine reduce the work of our algorithm by an asymptotic factor of 64. Thus, our MINFIT algorithm is an effective solver for least squares problems.

5.6. Data Structures.

Let us first assume that $n \leq 64$. As our algorithms access A by rows, we lay out the rows of the matrix across the processing elements of the ILLIAC. We thus represent the matrix by a sword vector $A[*]$ of order m .

We work with the rows of A to compute (a) the pairwise inner product, and (b) the new rows after a plane rotation. The GLYPNIR language provides a built-in function ROWSUM that sums the 64 numbers of a sword in 6 additions. The GLYPNIR expression

$$\text{ROWSUM}(A[I] * A[J])$$

computes the inner product of the i -th and j -th rows of A . If $n < 64$, we must disable the last $(64 - n)$ processing elements when we call the ROWSUM function. An alternative is to apply our algorithms to an $m \times 64$ matrix \hat{A} , given by

$$\hat{A} = (A|0) .$$

The following lines of GLYPNIR code compute the new i -th and j -th rows of A at the end of a plane rotation:

```
T      := A[I] * COSPHI + A[J] * SINPHI ;  
A[J]   := -A[I] * SINPHI + A[J] * COSPHI ;  
A[I]   := T ;
```

where T is a sword used for temporary storage.

We now consider the case when $n > 64$. Let

$$\ell = \left\lceil \frac{n}{64} \right\rceil ,$$

i.e. ℓ equals the smallest integer $\geq n/64$. We construct an $m \times 64\ell$ matrix \hat{A} , given by

$$\hat{A} = (A|0) .$$

The rows of \hat{A} are then divided into ℓ equal segments:

$$\hat{A} = (A_1|A_2|\cdots|A_\ell) .$$

Thus, we represent the matrix \hat{A} by the ℓ sword vectors $A1[*]$, \dots , $A\ell[*]$, each of order m .

The GLYPNIR expression

$$\text{ROWSUM}(A1[I] * A1[J] + \cdots + A\ell[I] * A\ell[J])$$

computes the inner product of the i -th and j -th rows of the matrix A . Plane rotations are applied to individual segments of the rows. For example, we write the ℓ lines of code

```
A1[J] := -A1[I] * SINPHI + A1[J] * COSPHI ;
A2[J] := -A2[I] * SINPHI + A2[J] * COSPHI ;
      ⋮
Aℓ[J] := -Aℓ[I] * SINPHI + Aℓ[J] * COSPHI ;
```

to compute the new vector

$$\tilde{a}_j := -\tilde{a}_i \cdot \sin \varphi + \tilde{a}_j \cdot \cos \varphi .$$

Since the columns of the matrix U are transformed in the same manner as the rows of A , we lay out U so that its columns lie across the processing elements. Thus, we represent the matrix by a sword vector $U[*]$ of dimension m . For the two different cases of $m \leq 64$ and $m > 64$, we apply techniques similar to those discussed in the previous paragraphs.

The rows of the data matrix B are modified in an identical fashion as the rows of A . Therefore, we lay out the rows of B across the processing elements. The two cases of column dimension $s \leq 64$ and $s > 64$ for B are dealt with in the same manner as the corresponding cases for the matrix A . We note that the execution time of Algorithm MINFIT is independent of s , for $s \leq 64$. If s is much greater than m , then the execution time will be proportional to $\left\lceil \frac{s}{64} \right\rceil$.

5.7. Numerical Properties

Let us examine the question of numerical stability. An error analysis of the action of plane rotations on a matrix was given by Wilkinson in his classic text [18]. His error bounds were later improved by Gentleman [5]. We use their results to study the effects of the plane rotations in one sweep of our algorithm.

Let

$$M = \frac{1}{2} m(m-1) , \quad (7.1)$$

and let R_j represent the j -th plane rotation, for $j = 1, 2, \dots, M$.

We can show that the computed matrix \tilde{A}_M after one sweep of rotations satisfies the inequality

$$\|\tilde{A}_M - R_M R_{M-1} \cdots R_1 A\| \leq 2^{-48} (m+n-2) (1 + 2^{-48})^{m+n-2} \|A\| . \quad (7.2)$$

The right-hand side of the inequality (7.2) is an extreme upper bound. We expect the statistical distribution of the rounding errors to reduce the error to well below the level of the bound; for this reason alone, a factor of the order of $(m+n-2)^{1/2}$ in place of $(m+n-2)$ might be more realistic. We see that our algorithm is extremely stable.

As the matrix U is formed as a product of plane rotations, we examine here the deviation from orthogonality of such a product. Let \tilde{Q}_M represent the computed product of the plane rotations in one sweep. We have the inequality that

$$\|\tilde{Q}_M - R_M R_{M-1} \cdots R_1\| \leq 2^{-48} m^{1/2} (m+n-2) (1+2^{-48})^{m+n-2}. \quad (7.3)$$

Again statistical consideration indicates that a factor of the order of $m^{1/4} (m+n-2)^{1/2}$ instead of $m^{1/2} (m+n-2)$ is probably more realistic. The matrix U is thus very close to an orthogonal matrix.

The tolerance τ controls the accuracy of the solution. At convergence of our algorithm, we have that

$$\|V_r^t V_r - I_r\| \leq r^{1/2} \tau. \quad (7.4)$$

Indeed, our numerical experiments show that the accuracy of the computed singular values and vectors of A is of the order of τ .

5.8. Test Results

We have written two GLYPNIR programs implementing Algorithms SVD and MINFIT. Tests were carried out on the ILLIAC IV computer.

EXAMPLE 1 (see [8]).

We have chosen the following matrices:

$$A = \begin{bmatrix} 22 & 10 & 2 & 3 & 7 \\ 14 & 7 & 10 & 0 & 8 \\ -1 & 13 & -1 & -11 & 3 \\ -3 & -2 & 13 & -2 & 4 \\ 9 & 8 & 1 & -2 & 4 \\ 9 & 1 & -7 & 5 & -1 \\ 2 & -6 & 6 & 5 & 1 \\ 4 & 5 & 0 & -2 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} -1 & 1 & 0 \\ 2 & -1 & 1 \\ 1 & 10 & 11 \\ 4 & 0 & 4 \\ 0 & -6 & -6 \\ -3 & 6 & 3 \\ 1 & 11 & 12 \\ 0 & -5 & -5 \end{bmatrix}.$$

The singular values of A are $\sqrt{1248}$, 20, $\sqrt{384}$, 0 and 0. Our SVD program computed those values to machine precision. The minimum norm solution to the overdetermined system

$$AX = B$$

is given by

$$X = \begin{pmatrix} -\frac{1}{12} & 0 & -\frac{1}{12} \\ 0 & 0 & 0 \\ \frac{1}{4} & 0 & \frac{1}{4} \\ -\frac{1}{12} & 0 & -\frac{1}{12} \\ \frac{1}{12} & 0 & \frac{1}{12} \end{pmatrix}.$$

Our MINFIT program returned a solution accurate to 14 decimal digits.

EXAMPLE 2 (see [6]).

We have chosen the matrix

$$A = \begin{pmatrix} 1 & -1 & -1 & -1 & -1 & \dots \\ & 1 & -1 & -1 & -1 & \dots \\ & & 1 & -1 & -1 & \dots \\ & & & 1 & -1 & \dots \\ & \bigcirc & & & 1 & \dots \\ & & & & & \dots \end{pmatrix}_{n \times n},$$

which is ill-conditioned as it has a very small singular value.

The matrix becomes singular if we add -2^{-n+2} to its $(n,1)$ position.

We applied our SVD program to this choice of A for different values of n . For comparison, we have chosen the SVD subroutine in the EISPACK eigenvalue package from the Argonne National Laboratory [3]. The EISPACK routine implements the method of Golub and Reinsch [8] and has been coded for high execution efficiency. We applied the routine to the matrix A on an IBM 370/168 computer at the Stanford Linear Accelerator Center. The code was compiled by the FORTRAN H EXTENDED compiler with optimization level 2.

n	ILLIAC IV		IBM 370/168	ILLIAC TIME
	iter	time	time	IBM TIME
16	7	0.26	0.101	2.57
32	8	1.25	0.57	2.19
48	8	2.89	1.76	1.64
64	9	5.57	4.03	1.38
96	10	15.94	12.81	1.24
128	9	26.81	29.68	0.90

We should point out that the GLYPNIR compiler produces very inefficient code. We have written another program in CFD implementing our SVD algorithm. The ILLIAC execution time with $n = 64$ was 3.31 seconds, a saving of 41% over the GLYPNIR code.

It is unfortunate that due to certain limitations we were not able to run examples with larger values of n on the ILLIAC. Nonetheless, we observe that our ILLIAC routine becomes more efficient relative to the EISPACK routine with increasing values of n . The execution time of the former is crudely proportional to

$$\left(\text{iter} \times \left\lceil \frac{n}{64} \right\rceil \times n^2 \right),$$

and that of the latter to n^3 . There is thus a great potential in matrix computations of a parallel computer with many processors.

REFERENCES

- [1] Chan, T. F. C., "On computing the singular value decomposition," Report STAN-CS-77-588, Computer Science Dept., Stanford University (1977).
- [2] Chartres, B.A., "Adaptation of the Jacobi method for a computer with a magnetic-tape backing store," Computer J. 5 (1962), 51-60.
- [3] Garbow, B.S., Boyle, J.M., Dongarra, J.J., and Moler, C.B., Matrix Eigensystem Routines--EISPACK Guide Extension Springer-Verlag, Berlin (1977).
- [4] Gentleman, W.M., "Least squares computations by Givens' transformations without square roots," J. Inst. Maths. Applics. 12 (1973), 329-336.
- [5] Gentleman, W.M., "Error analysis of QR decomposition by Givens transformations," Lin. Alg. Applics. 10 (1975), 189-197.
- [6] Golub, G.H., and Kahan, W., "Calculating the singular values and pseudoinverse of a matrix," J. SIAM Ser. B: Numer. Anal. 2 (1965), 205-224.
- [7] Golub, G.H., and Luk, F.T., "Singular value decomposition: applications and computations," ARO Report 77-1, Transactions of the 22-nd Conference of Army Mathematicians (1977), 577-605.
- [8] Golub, G.H., and Reinsch, C., "Singular value decomposition and least squares solutions," Numer. Math. 14 (1970), 403-420.
- [9] Heller, D., "A survey of parallel algorithms in numerical linear algebra," Report, Computer Science Dept., Carnegie-Mellon University (February 1976).
- [10] Hestenes, M.R., "Inversion of matrices by biorthogonalization and related results," J. Soc. Indust. Appl. Math. 6 (1958), 51-90.
- [11] Lawrie, D.H., Layman, T., Baer, D., and Randal, J.M., "GLYPNIR--a programming language for ILLIAC IV," Comm. ACM 18 (1975), 157-164.
- [12] Nash, J.C., "A one-sided transformation method for the singular value decomposition and algebraic eigenproblem," Computer J. 18 (1975), 74-76.
- [13] Rutishauser, H., "The Jacobi method for real symmetric matrices," Numer. Math. 2 (1966), 1-10.

- [14] Sameh, A.H., and Kuck, D.J., "A parallel QR algorithm for tri-diagonal symmetric matrices," Report, Computer Science Dept., University of Illinois, Urbana-Champaign (July 1974).
- [15] Stevens, K.G., Jr., "CFD--a FORTRAN-like language for the ILLIAC IV," ACM Sigplan Notices 10 (1975), 72-76.
- [16] Stewart, G.W., private communication (1978).
- [17] Wilkinson, J.H., "Note on the quadratic convergence of the cyclic Jacobi process," Numer. Math. 4 (1962), 296-300.
- [18] Wilkinson, J.H., The Algebraic Eigenvalue Problem, Clarendon, Oxford (1965).